

# VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ

BRNO UNIVERSITY OF TECHNOLOGY

FAKULTA INFORMAČNÍCH TECHNOLOGIÍ  
ÚSTAV POČÍTAČOVÉ GRAFIKY A MULTIMÉDIÍ

FACULTY OF INFORMATION TECHNOLOGY  
DEPARTMENT OF COMPUTER GRAPHICS AND MULTIMEDIA

## ROZPOZNÁVAČ ŘEČI ŘÍZENÝ GRAMATIKAMI

DIPLOMOVÁ PRÁCE  
MASTER'S THESIS

AUTOR PRÁCE  
AUTHOR

Bc. VOJTĚCH ŠKORVAGA

BRNO 2014



**VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ**  
BRNO UNIVERSITY OF TECHNOLOGY



**FAKULTA INFORMAČNÍCH TECHNOLOGIÍ**  
**ÚSTAV POČÍTAČOVÉ GRAFIKY A MULTIMÉDIÍ**

FACULTY OF INFORMATION TECHNOLOGY  
DEPARTMENT OF COMPUTER GRAPHICS AND MULTIMEDIA

## **ROZPOZNÁVAČ ŘEČI ŘÍZENÝ GRAMATIKAMI**

GRAMMAR BASED AUTOMATIC SPEECH RECOGNIZER

### **DIPLOMOVÁ PRÁCE**

MASTER'S THESIS

### **AUTOR PRÁCE**

AUTHOR

**Bc. VOJTĚCH ŠKORVAGA**

### **VEDOUCÍ PRÁCE**

SUPERVISOR

**Ing. PETR SCHWARZ, Ph.D.**

BRNO 2014

**Vysoké učení technické v Brně - Fakulta informačních technologií**

Ústav počítačové grafiky a multimédií

Akademický rok 2013/2014

**Zadání diplomové práce**

Řešitel: **Škorvaga Vojtěch, Bc.**

Obor: Management a informační technologie

Téma: **Rozpoznávač řeči řízený gramatikami**

**Grammar Based Automatic Speech Recognizer**

Kategorie: Zpracování řeči a přirozeného jazyka

Pokyny:

1. Prostudujte formát gramatiky SGRS, rozpoznávač firmy Phonexia (aplikační rozhraní a formát rozpoznávacích sítí), způsob generování rozpoznávacích sítí pomocí konečných převodníků, síťový protokol MRCP v2 a open source projekt UniMRCP.
2. Napište algoritmus na převod SGRS na konečný automat.
3. Dle příkladu naimplementujte algoritmus na složení rozpoznávací sítě na základě konečného automatu s gramatikou, převodníku popisujícího lexikon, převod na kontextově závislé fonémy, a na stavy.
4. Otestujte kompilátor gramatik v součinnosti s rozpoznávačem.
5. Napište modul pro projekt UniMRCP, který zapouzdřuje rozpoznávač řeči a přidává komunikaci s rozpoznávačem pomocí síťového protokolu MRCP v2.
6. Otestujte celkové řešení za pomoci některé IVR platformy (Asterisk, FreeSwitch, OptimSys)

Literatura:

- Dle konzultace se školitelem

Při obhajobě semestrální části diplomového projektu je požadováno:

- Body 1 a 2.

Podrobné závazné pokyny pro vypracování diplomové práce naleznete na adrese

<http://www.fit.vutbr.cz/info/szz/>

Technická zpráva diplomové práce musí obsahovat formulaci cíle, charakteristiku současného stavu, teoretická a odborná východiska řešených problémů a specifikaci etap, které byly vyřešeny v rámci ročníkového a semestrálního projektu (30 až 40% celkového rozsahu technické zprávy).

Student odevzdá v jednom výtisku technickou zprávu a v elektronické podobě zdrojový text technické zprávy, úplnou programovou dokumentaci a zdrojové texty programů. Informace v elektronické podobě budou uloženy na standardním nepřepisovatelném paměťovém médiu (CD-R, DVD-R, apod.), které bude vloženo do písemné zprávy tak, aby nemohlo dojít k jeho ztrátě při běžné manipulaci.

Vedoucí: **Schwarz Petr, Ing., Ph.D.,** UPGM FIT VUT

Datum zadání: 1. listopadu 2013

Datum odevzdání: 28. května 2014

**VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ**  
Fakulta informačních technologií  
Ústav počítačové grafiky a multimédií  
602 00 Brno, Božetěchova 2



doc. Dr. Ing. Jan Černocký  
vedoucí ústavu

## Abstrakt

Tato práce popisuje vytvoření systému pro sestavení rozpoznávací sítě pro rozpoznávač mluvené řeči na základě Speech Recognition Grammar Specification (SRGS) gramatiky definované W3C konsorciem. Rozpoznávač byl spolu s tímto modulem integrován do softwarové ústředny FreeSWitch pomocí kombinace síťových protokolů MRCPv2/SIP/RTP a testován. Práce byla řešená ve spolupráci s firmou Phonexia s.r.o.

## Abstract

This work describes a development of system for network compilation for speech recognition based on Speech Recognition Grammar Specification (SRGS) grammar defined by W3C consortium. Together with the new module, the recognizer was integrated to the FreeSwitch software phone switch using a combination of MRCPv2/SIP/RTP networks protocols and tested.

## Klíčová slova

SRGS, OpenFST, Phonexia s.r.o., UniMRCP, MRCP, IVR, FreeSwitch

## Keywords

SRGS, OpenFST, Phonexia s.r.o., UniMRCP, MRCP, IVR, FreeSwitch

## Citace

Vojtěch Škorvaga: Rozpoznávač řeči řízený gramatikami, diplomová práce, Brno, FIT VUT v Brně, 2014

# Rozpoznávač řeči řízený gramatikami

## Prohlášení

Prohlašuji, že jsem tuto práci k semestrálnímu projektu vypracoval samostatně pod vedením pana Ing. Petra Schwarze, Ph.D.

.....

Vojtěch Škorvaga  
28. května 2014

© Vojtěch Škorvaga, 2014.

*Tato práce vznikla jako školní dílo na Vysokém učení technickém v Brně, Fakultě informačních technologií. Práce je chráněna autorským zákonem a její užití bez udělení oprávnění autorem je nezákonné, s výjimkou zákonem definovaných případů.*

# Obsah

<b>1 Úvod</b>	<b>4</b>
<b>2 Teoretická východiska práce</b>	<b>5</b>
2.1 Phonexia	5
2.1.1 Tvorba rozpoznávacích sítí	6
2.2 Rozpoznávač řeči	6
2.2.1 Dekodér	6
2.3 Rozpoznávací síť	7
2.3.1 Definice konečného automatu	8
2.3.2 Definice váhovaný konečného automatu	8
2.3.3 Definice váhovaného převodníku	8
2.3.4 Vztah váhovaného konečného automatu a převodníků	8
2.3.5 Komponenty rozpoznávací sítě	9
2.3.6 Skládání rozpoznávací sítě	11
2.4 Gramatika pro rozpoznávání řeči	13
2.4.1 ABNF forma	14
2.4.2 XML forma	14
2.4.3 Elementy SRGS	16
2.5 Hlasový dialogový systém	16
2.5.1 Obecné schéma IVR systému	17
2.5.2 Existující řešení	17
2.6 Protokol MRCP	17
2.6.1 Proces rozpoznávání v MRCP	18
2.6.2 UniMRCP	18
<b>3 Převod SGRS gramatik na váhovaný konečný automat</b>	<b>20</b>
3.1 Vstup	22
3.2 Výstup	22
3.3 Postup převodu gramatiky na graf G	22
3.4 ANTRL	23
3.4.1 ANTRL pravidla pro SGRS/ABNF	24
3.4.2 Zpracování chyb	26
<b>4 Sestavení rozpoznávací sítě</b>	<b>28</b>
4.1 Formát rozpoznávací sítě v BSAPI	28
4.1.1 Převod na akceptor <i>HCLGa</i>	29
4.1.2 Převod na formát pro BSAPI	30
4.1.3 Optimalizace BSAPI sítě	32

4.1.4	Vytvoření převodníku L . . . . .	33
4.2	Zapouzdření do BSAPI . . . . .	33
<b>5</b>	<b>Součinnost gramatik s rozpoznávačem</b>	<b>36</b>
5.1	Testy sestavení rozpoznávací sítě . . . . .	36
5.1.1	Rychlost načtení gramatiky a sestavení rozpoznávací sítě . . . . .	36
5.1.2	Přesnost rozpoznávací sítě . . . . .	37
5.1.3	Rychlost rozpoznávací sítě . . . . .	37
5.1.4	Paměťové nároky rozpoznávací sítě . . . . .	38
5.1.5	Zhodnocení testů . . . . .	39
<b>6</b>	<b>Rozpoznávač jako modul UniMRCP</b>	<b>40</b>
6.1	Plugin UniMRCP . . . . .	40
6.1.1	Vytvoření nového pluginu . . . . .	40
6.1.2	Konfigurace pluginu . . . . .	40
6.1.3	Rozhraní pluginu rozpoznávače . . . . .	41
6.1.4	Ukončení rozpoznávání . . . . .	41
<b>7</b>	<b>Použití UniMRCP ve FreeSwitch</b>	<b>43</b>
7.1	Konfigurace spojení s UniMRCP . . . . .	43
7.2	Dialogový scénář . . . . .	44
7.3	Přiřazení dialogového scénáře k telefonnímu číslu . . . . .	44
<b>8</b>	<b>Závěr</b>	<b>46</b>
<b>A</b>	<b>Obsah CD</b>	<b>48</b>

# Seznam obrázků

2.1	Schéma rozpoznávače řeči. . . . .	7
2.2	Vizualizace rozpoznávací sítě pro anglické slovo "ONE" nebo "ZERO" v OpenFST. . . . .	7
2.3	Vizualizace konečného automatu G, reprezentující gramatiku v OpenFST. . . . .	9
2.4	Vizualizace převodníku L pro anglické číslovky v OpenFST. . . . .	10
2.5	Vizualizace převodníku C v OpenFST. . . . .	11
2.6	Vizualizace výřezu z převodníku H v OpenFST. . . . .	12
2.7	Příklad skládání převodníků. (a) převodník A, (b) převodník B, (ab) složený převodník $A \circ B$ . Obrázky jsou přejaty z projektu OpenFST. . . . .	13
2.8	Příklad determinizace převodníků. (c) převodník C, (d) determinizovaný převodník C. Obrázky jsou přejaty z projektu OpenFST. . . . .	13
2.9	Příklad minimalizace převodníků. (e) převodník E, (f) minimalizovaný převodník E. Obrázky jsou přejaty z projektu OpenFST. . . . .	14
2.10	Schéma IVR systému. . . . .	17
2.11	Stavový automat rozpoznávání. . . . .	18
3.1	Schéma elementu slovo v VKA G. . . . .	20
3.2	Schéma elementu pro definování pravidla v VKA G. . . . .	20
3.3	Schéma elementu pro výběr alternativy v VKA G. . . . .	21
3.4	Schéma elementu pro vybrání jednu z možností nebo nic v VKA G. . . . .	21
3.5	Schéma elementu pro vybrání jednu z možností nebo nic v VKA G. . . . .	21
3.6	Schéma elementu pro nekonečné opakování v VKA G. . . . .	22
3.7	Schéma elementu pro začátek gramatiky v VKA G. . . . .	22
3.8	Vizualizace převodníku G (zpracované gramatiky) v OpenFST. . . . .	23
4.1	Zapozdření tříd pro sestavení rozpoznávací sítě do BSAPI. . . . .	35
6.1	Diagram UniMRCP modulu pro ASR. . . . .	42



# Kapitola 1

## Úvod

Řeč je prostředek k interakci mezi lidmi již stovky tisíc let. Pomohla lidstvu sdělovat myšlenky, tvořit jazyky, kultury i národy. Přestože obraz vydá za tisíce slov, je to řeč, kterou ho můžeme popsat jinému člověku. A nyní, mohou lidé řečí komunikovat i se stroji. Stroje mohly mluvit na člověka již delší dobu, třeba díky jednoduchému pouštění audio záznamu. Opačný směr až v poslední době, díky růstu výpočetního výkonu a pokroku v analýze zpracování přirozeného jazyka.

V tomto textu se budu věnovat rozšíření rozpoznávače řeči od firmy Phonexia s.r.o. o prostředky o možnost definice rozpoznávaných vět pomocí W3C Speech Recognition Grammar Specification (SRGS) gramatik následně bude rozpoznávač a nový UniMRCP modul integrován do softwarové ústředny FreeSwitch pomocí protokolů MRCPv2/SIP/RTP. pro použití rozpoznávače řeči, který je součástí BSAPI, v řečových automatech postavených nad open source projekty jako UniMRCP a FreeSWITCH.

Cílem projektu je vytvořit funkční dialogový systém s možností popisu rozpoznávaných vět pomocí SRGS gramatiky. K tomu zapotřebí vytvořit převodník SRGS gramatik na konečný automat, kompilaci rozpoznávací sítě a vytvořit plugin pro UniMRCP projekt s využitím rozpoznávače firmy Phonexia.

## Kapitola 2

# Teoretická východiska práce

V této kapitole je popsán vztah SRGS (Speech Recognition Grammar Specification) k rozpoznávání řeči, hlasový dialogový systém (Interactive Voice Response, IVR) a role rozpoznávače řeči (Automatic Speech Recognition, ASR) v tomto systému. Jde o obecné předpoklady pro vytvoření IVR systému s ASR řízeného pomocí SRGS gramatik pro společnost Phonexia.

### 2.1 Phonexia

Vyvíjí technologie pro automatické získávání informací z řeči (identifikace jazyka, řečníka, detekce klíčových slov a přepis řeči na text). Produkty se uplatňují při analýze telefonátů v call-centrech a v oblasti bezpečnosti nebo obrany.

Firma vyvíjí SDK (Software Development Kit) pro své produkty v C/C++, které pojmenovala BSAPI(Brno Speech API). Jeho součástí jsou technologie jako:

**VAD** (Voice Activity Detector) Systém pro detekci lidského hlasu v audiu.

**SID** (Speaker Identification) Systém pro identifikaci řečníka.

**LID** (Language Identification) Systém identifikace jazyka z audia.

**KWS** (Key Word Spotting) Systém pro vyhledávání klíčových slov v audiu.

**STT** (Speech To Text) Systém pro přepis řeči na text, také známý jako.

Tyto technologie mohou pracovat offline nebo online. Offline znamená, že se zpracovávají soubory s pořízenými nahrávkami, zatímco online zpracování znamená, že program čte data přímo z mikrofonu a pro analýzu nemůže využít zbytek promluvy, která ještě nebyla vyřčena.

Vstupem pro KWS a STT (souhrně ASR, Automatic Speech Recognition) je takzvaná rozpoznávací síť, která definuje co se bude hledat. V technologii KWS může tato síť obsahovat pouze hledaná klíčová slova, zatímco v STT obsahuje několik set tisíc až miliónu slov ze slovníku pro zvolený jazyk a popis jejich vazeb. Pro technologii KWS si tuto síť nadefinuje uživatel.

Rozpoznávací síť je datová struktur popisující konečný automat.

### 2.1.1 Tvorba rozpoznávacích sítí

Se označuje za kompilaci. Na začátku projektu nebyla tato funkce součástí rozpoznávacího jádra (BSAPI). Tuto funkcionalitu zajišťovala sada externích skriptů pro linuxovou příkazovou řádku, s využitím několika programů, pro převod různých formátů a pro manipulaci s konečnými automaty. Skripty byly dodávány zákazníkům, ale jsou značně omezující.

Hlavní nevýhodou skriptové kompilace sítí je její platformní závislost. Linuxové skripty budou těžko fungovat na Windows. Díky tomu je nelze integrovat do systémů zákazníka. Dále je velmi omezující proces kompilace z uživatelského hlediska, probíhá následovně:

1. Spuštění.
2. Načtení slovníků spojení, které má síť rozpoznávat.
3. Načtení akustických modelů pro zvolený jazyk. Velikost souboru s modelem se pohybuje kolem 200 MB.
4. Kompilace sítí pomocí OpenFST<sup>1</sup>. Tato část je inspirovaná projektem KALDI<sup>2</sup> jehož popis je v článku [4].
5. Uložení sítě do souboru.
6. Načtení sítě ze souboru do BSAPI.

Tento proces se opakuje pro každou uživatelem navrženou síť. Data se mezi jednotlivými fázemi kompilace předávají pomocí souborů. Proto je proces kompilace značně pomalý a těžkopádný.

K jeho ergonomii nepřispívá ani metoda zadávání slovníkových spojení, které má rozpoznávač rozpoznávat. Její podstatou je výčet všech možností. Při takovém způsobu zápisu by měl vstupní soubor pro rozpoznání hesla, které může mít 12 číslic, asi  $1e^{12}$  řádků.

Z těchto důvodů je cílem této diplomové práce vytvořit v BSAPI kompilátor rozpoznávacích sítí, který jako vstup bude akceptovat gramatiku z kapitoly definovanou W3C SRGS.

## 2.2 Rozpoznávač řeči

Zachycuje obrázek 2.1. Z řeči jsou extrahovány příznaky. K těm se potom hledá odpovídající index do tabulky, která tvoří akustický model. Tyto informace se posílají do dekodéru, který předtím načte rozpoznávací síť.

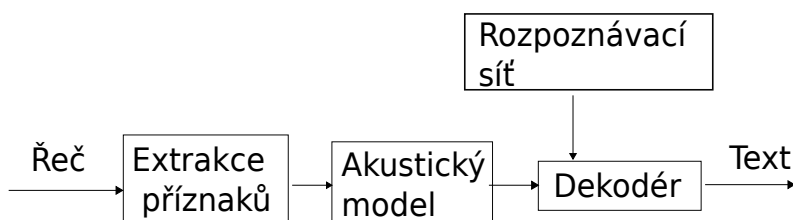
### 2.2.1 Dekodér

Pracuje, tak že na základě rozpoznávací sítě a výstupu z porovnání akustického modelu s řečí tvoří hypotézy. Které dále zpřesňuje a eliminuje nepravděpodobné průchody rozpoznávací sítě.

Samotný text dekodér získá tak že z porovnání akustického modelu s řečí získá indexy do tabulky akustického modelu. Tyto indexy reprezentují kontextově závislé fonémy, které se

<sup>1</sup>OpenFST je knihovna pro reprezentaci konečných automatů. Domovská stránka projektu je na webu <http://www.openfst.org/twiki/bin/view/FST/WebHome>

<sup>2</sup>Představení toolkitu na webu: <http://www.fit.vutbr.cz/units/UPGM/prod/index.php?id=304&notitle=1>



Obrázek 2.1: Schéma rozpoznávače řeči.

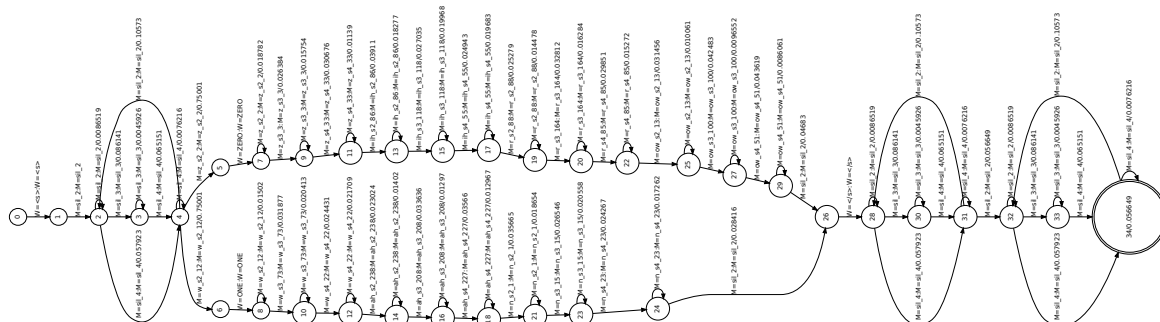
převodou na kontextově nezávislé. A posloupnost kontextově nezávislých fonémů odpovídá slovu.

Dekodér pracuje s více možnými výsledky. Jako výsledek prohlásí ten s nejvyšší pravděpodobností.

## 2.3 Rozpoznávací síť

Se podle článku [7] skládá ze 3 převodníků (převodník vysvětlen v 2.3.3) a jednoho konečného automatu. Převodník<sup>3</sup> je graf který na svých hranách rozlišuje vstupní a výstupní řetězec. Průchod automatem do konečného stavu pak mapuje jednu větu na druhou. Jinak řečeno převede jeden vzor na druhý od toho název převodník.

Na obrázku 2.2 je rozpoznávací síť, která popisuje promluvu jednoho z anglických slov "one" nebo "zero". Která byla vytvořena podle článku [7]. Síť obsahuje akustické



Obrázek 2.2: Vizualizace rozpoznávací sítě pro anglické slovo "ONE" nebo "ZERO" v OpenFST.

jednotky a slova. Akustické jednotky jsou na vstupech hran převodníků a slova na jejich výstupu.

<sup>3</sup>Anglicky transducer.

### 2.3.1 Definice konečného automatu

Zkráceně KA, je pětice:

$KA = (Q, \Sigma, \sigma, s, F)$ , kde:

$Q$  je konečná množina vnitřních stavů automatu.

$\Sigma$  konečná množina vstupních symbolů automatu, také nazývaná abeceda.

$\sigma$  je přechodovou funkcí ze stavu do stavu, definovanou jako  $Q \times \{\Sigma \cup \epsilon\} \rightarrow 2^Q$

$s$  je množina počátečních stavů, taková že  $s \subseteq Q$

$F$  je množina koncových stavů, taková že  $F \subseteq Q$

### 2.3.2 Definice váhovaný konečného automatu

Akceptor je vlastně konečný automat, ale pro potřeby této práce je ho musíme rozšířit o váhu z okruhu  $K$ . Okruh reprezentuje váhu s jakou se může uskutečnit přechod nebo akceptovat koncový stav.

Upravená definice tedy vypadá jako uspořádaná šestice:  $A = (Q, \Sigma, \sigma, s, F, p)$ , kde:

$Q$  je konečná množina vnitřních stavů automatu.

$\Sigma$  konečná množina vstupních symbolů automatu, také nazývaná abeceda.

$\sigma$  je množina přechodů, ve tvaru  $Q \times \{\Sigma \cup \epsilon\} \times K \rightarrow 2^Q$ .  $s$  je množina počátečních stavů, taková že  $s \subseteq Q$

$F$  je množina koncových stavů, taková že  $F \subseteq Q$

$p$  je váhová funkce koncových stavů, taková že  $p : F \rightarrow K$

### 2.3.3 Definice váhovaného převodníku

Váhovaný Převodník je váhovaný konečný automat, rozšířený o výstupní symbol na hraně. Akceptováním nějaké věty pak dojde k jejímu přepisu z jazyka definovaného konečným automatem a jeho vstupními symboly do jazyka definovaného konečným automatem a jeho výstupními symboly.

V článku [7] a [6] je převodník s váhami nad okruhem  $K$  definován jako osmice:  $T = (\Sigma, \Delta, Q, I, F, E, \lambda, p)$ , kde

$\Sigma$  je konečná abeceda vstupních symbolů.

$\Delta$  je konečná abeceda výstupních symbolů.

$Q$  je konečná množina vnitřních stavů automatu.

$I$  je množina počátečních stavů, taková že  $I \subseteq Q$

$F$  je množina koncových stavů, taková že  $I \subseteq Q$

$E$  je množina přechodů, taková že  $E \subseteq Q \times (\Sigma \cup \{\epsilon\}) \times (\Delta \cup \{\epsilon\}) \times K \times Q$

$\lambda$  je váhová funkce počátečních stavů, taková že  $\lambda : I \rightarrow K$

$p$  je váhová funkce koncových stavů, taková že  $p : F \rightarrow K$

### 2.3.4 Vztah váhovaného konečného automatu a převodníků

Váhovaný konečný automat je současně speciálním typem váhovaného převodníku, který má na vstupní a výstupní hraně stejný symbol. Váhovaný konečný automat tedy lze jednoduše převést na převodník:

**Data:** Váhováný konečný automat  $A = (Q_A, \Sigma, \sigma, s, F_A, p_A)$   
**Result:** Váhováný převodník  $T = (\Sigma, \Sigma, Q_A, s, F_A, E, (s, 1), p_a)$   
method  
**foreach**  $(start_s, symbol, likelihood, end_s) \in \sigma$  **do**  
|  $E \leftarrow E \cup \{(start_s, symbol, symbol, likelihood, end_s)\}$   
**end**

**Algorithm 1:** Převod váhovaného KA na váhovaný převodník.

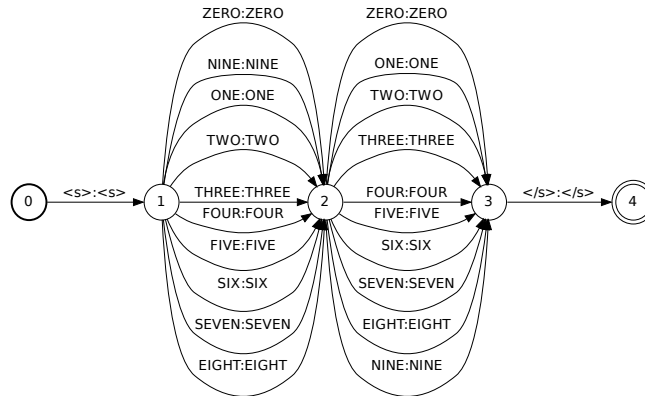
### 2.3.5 Komponenty rozpoznávací sítě

Při skládání rozpoznávací sítě se podle článku [7] využívá těchto váhovaných převodníků:

- G konečný automat popisující gramatiku. Tento graf je výstupem z analýzy SRGS gramatiky.
- L lexikální převodník. Mapuje fonémy na slova. Fonémy jsou reprezentované třeba fonetickou abecedou.
- C kontextový převodník. Mapuje fonémy kontextově fonémy.
- H převodník mapuje kontextově závislé fonémy na stavy. Které jsou popsány indexy do tabulky stavů skrytých Markovových modelů (HMM<sup>4</sup>)

#### Komponenta G

Je váhovaný konečným automatem reprezentující uživatelem zadanou SRGS gramatiku. Obrázek 2.3 ukazuje konečný automat pro gramatiku přijímající dvě číslovky. Na hranách automatu jsou vidět čísla definovaná uživatelem.

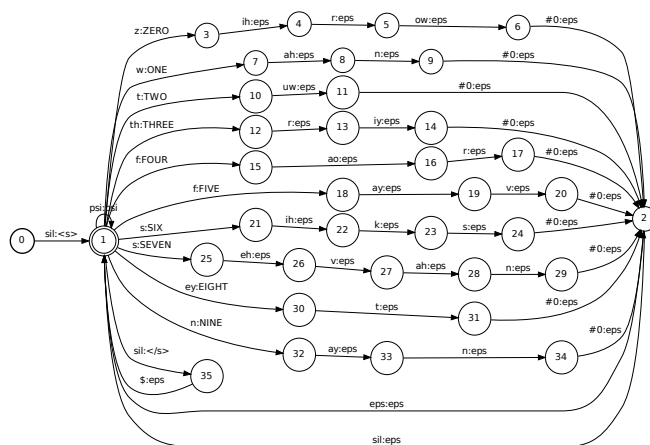


Obrázek 2.3: Vizualizace konečného automatu G, reprezentující gramatiku v OpenFST.

<sup>4</sup>Anglicky Hidden Markov Model.

## Komponenta L

Reprezentuje výslovnost slov v konečném automatu G (gramatice), tím že převádí fonémy na slova. Na obrázku 2.4 jsou vidět fonémy na vstupech hran a slovo nebo symbol pro  $\varepsilon$  na výstupech hrany. Protože vychází z výslovnosti tak jedno slovo může mít více výslovností a tím pádem i více fonetických prepisů<sup>5</sup>. Nebo naopak rozdílná slova s jedním fonetickým prepisem, která jsou v češtině označována za homonyma. V případě přidání fonetického prepisu se problém řeší jednoduše pouhým vložením další větve ze stavu 1 do stavu 2. V případě homonym je to složitější. Aby bylo možno převodníkem pracovat, tak článek v [7] autor přidává na konec každého fonetického prepisu homonym pomocný symbol. Tento symbol nemají žádné dva stejné znějící homonyma stejný. Symbol začíná znakem # a pokračuje číslem. V pozdějších fázích sestavení rozpoznávací sítě se odstraňuje převedením na symbol  $\varepsilon$ , aby nepřekážel dekodéru.



Obrázek 2.4: Vizualizace převodníku L pro anglické číslovky v OpenFST.

## Komponenta C

Váhouvaný převodník mapuje kontextově závislé fonémy na nezávislé fonémy. Využití kontextových fonémů je motivováno zvýšením přesnosti rozpoznávače řeči. Převodník má na vstupu hrany kontextově závislý foném a na výstupu nezávislý foném, které na vstupu hranách převodníků L. Příklad převodníků C pro tři fonémy je vidět na obrázku 2.5.

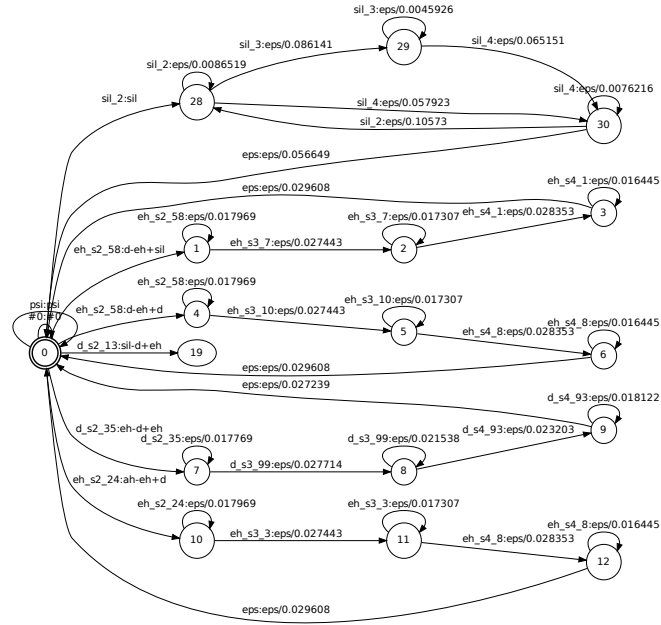
## Komponenta H

Váhouvaný převodník mapuje index do tabulky obsahující akustický model na kontextově závislé fonémy. Obrázek 2.6 zachycuje výřez převodníku H, který byl vytvořen pro kontextově závislé fonémy z obrázku 2.5. Na vstupu hran je vidět index a na výstupu je kontextově závislý foném.

<sup>5</sup>Fonetický prepis je prostý prepis slova na posloupnost fonémů kterými je vyslovován.







Obrázek 2.6: Vizualizace výřezu z převodníku H v OpenFST.

**Data:**  $T_1 = (A, B_1, Q_1, I_1, F_1, E_1, \lambda_1, p_1)$

$T_2 = (B_2, C, Q_2, I_2, F_2, E_2, \lambda_2, p_2)$

**Result:**  $T = (A, C, Q, I, F, E, \lambda, p)$

method

$Q \leftarrow I_1 \times I_2$

$S \leftarrow I_1 \times I_2$

**while**  $S \neq \emptyset$  **do**

$(q_1, q_2) \leftarrow \text{HEAD}(S)$

$\text{DEQUEUE}(S)$

**if**  $(q_1, q_2) \in I_1 \times I_2$  **then**

$I \leftarrow I \cup \{(q_1, q_2)\}$

$\lambda(q_1, q_2) \leftarrow \lambda_1(q_1) \otimes \lambda_2(q_2)$

**end**

**if**  $(q_1, q_2) \in F_1 \times F_2$  **then**

$F \leftarrow F \cup \{(q_1, q_2)\}$

$p(q_1, q_2) \leftarrow p_1(q_1) \otimes p_2(q_2)$

**end**

**for** each  $(e_1, e_2) \in E[q_1] \times E[q_2]$  such that  $B_1[e_1] = B_2[e_2]$  **do**

**if**  $(n[e_1], n[e_2]) \notin Q$  **then**

$Q \leftarrow Q \cup \{(n[e_1], n[e_2])\}$

$\text{ENDEQUEUE}(S, \{(n[e_1], n[e_2])\})$

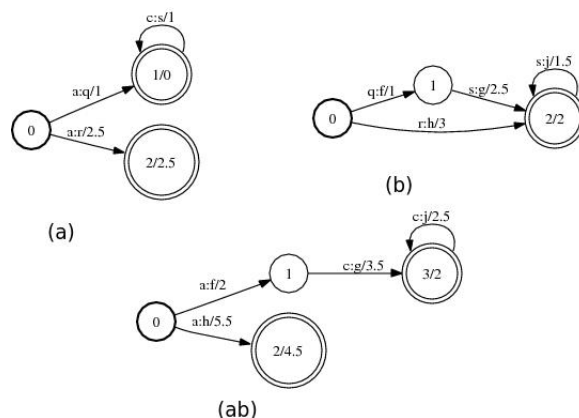
**end**

$E \leftarrow E \cup \{((q_1, q_2), A[e_1], C[e_2], \text{weight}(e_1) \otimes \text{weight}(e_2), (n[e_1], n[e_2]))\}$

**end**

**end**

**Algorithm 2:** Pseudo algoritmus pro skládání dvou převodníků.

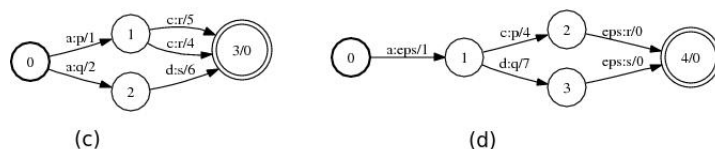


Obrázek 2.7: Příklad skládání převodníků. (a) převodník A, (b) převodník B, (ab) složený převodník  $A \circ B$ . Obrázky jsou přejaty z projektu OpenFST.

### Determinizace převodníku

Podle článku [7] a [6] obsahuje jeden stav maximálně jednu hranu s každým vstupním symbolem ze vstupní abecedy a  $\varepsilon$  se nevyskytuje na vstupu žádné hrany. Obrázek 2.8 ukazuje příklad determinizace. Determinizace způsobí, že pro každý vstupní řetězec existuje jenom jedna možná cesta automatem. Tím se urychlí analýza řeči v ASR.

Determinizovaný váhovaný převodník musí dávat pro stejný vstupní řetězec stejný výstup se stejnou pravděpodobností. Ne váhovaný převodník, nebo váhovaný konečný automat může být determinizován.



Obrázek 2.8: Příklad determinizace převodníků. (c) převodník C, (d) determinizovaný převodník C. Obrázky jsou přejaty z projektu OpenFST.

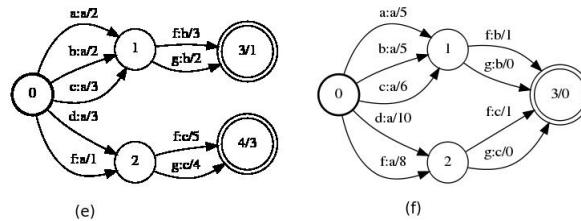
### Minimalizace váhovaného převodníku

Se provádí aby byla zmenšena jeho velikost a tím byl ušetřen čas během analýzy řeči v ASR. Na obrázku 2.9 je vidět příklad minimalizace.

## 2.4 Gramatika pro rozpoznávání řeči

Anglicky Speech Recognition Grammar, zkratka SRG. Někdy označována jako SRGS, kde poslední S je pro anglické slovo Specification. Je specifikace sady pravidel pro informování řečového rozpoznávače o tom, co čekáme, že uživatel řekne. Standard vytvořila a udržuje skupina W3C(World Wide Web Consortium) <sup>6</sup>.

<sup>6</sup>W3C je mezinárodní uskupení pro vývoj webových standardů, více na webu [www.w3.org](http://www.w3.org).



Obrázek 2.9: Příklad minimalizace převodníků. (e) převodník E, (f) minimalizovaný převodník E. Obrázky jsou přejaty z projektu OpenFST.

SRGS popisuje bezkontextovou gramatiku. Gramatika definovaná s pomocí SRGS je regulární gramatikou. Rozpoznávač je schopen pracovat pouze s regulární gramatikou definovanou konečným automatem na jehož hranách jsou slova a akustické jednotky, nebo  $\varepsilon$  pro reprezentaci žádného slova. Tento automat je uložen v rozpoznávací síti.

která je ekvivalentní s konečným automatem, na jehož hranách jsou očekávána slova nebo epsilon pro reprezentaci žádného slova.

Výstupem z rozpoznávače je pak řetězec slov odpovídající průchodu konečným automatem.

SRGS dovoluje gramatiky zapsat ve dvou ekvivalentních syntaxích:

- **ABNF** Rozšířena, anglicky Augmented, Backus/Naurova Forma<sup>7</sup>.
- **XML** forma.

#### 2.4.1 ABNF forma

Je syntaxe v čistém textu. ABNF je na rozdíl od XML zápis přehlednější a kratší. Díky tomu se s ním, novým uživatelům, lépe pracuje. Soubor v tomto formátu gramatiky by podle SRGS měl mít na svém začátku text "#ABNF" (bez uvozovek).

Příklad gramatiky pro objednání pizzy s pitím:

```
#ABNF 1.0 ISO-8859-1;
language en-US;
// root role is start point of grammar
root $pizza;
public $digit = one | two | three | four | five | six | seven | eight | nine;
public $number = $digit <1-> ;
public $kind = mushroom | salami | cheese | feferony;
public $drink = (big | small) colla | juice;
public $pizza = [ $number ] $kind [ pizza ] and [ $drink ];
```

#### 2.4.2 XML forma

Existuje protože XML je obecným standardem pro předávání dat. Jednotlivé prvky SRG jsou zapsány v XML elementech a jejich atributech. Soubor s touto formou gramatiky by podle SRGS měl mít na svém začátku "#XML" (bez uvozovek).

<sup>7</sup>ABNF je rozdílná od EBNF (Extended Backus/Naurova form). ABNF je definováno v RFC 5234 a EBNF v ISO 14977.

```

<?xml version="1.0" encoding="ISO-8859-1"?>
<!DOCTYPE grammar PUBLIC "-//W3C//DTD GRAMMAR 1.0//EN"
    "http://www.w3.org/TR/speech-grammar/grammar
        .dtd">

<grammar xmlns="http://www.w3.org/2001/06/grammar"
    xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
    xsi:schemaLocation="http://www.w3.org/2001/06/grammar
        http://www.w3.org/TR/speech-grammar/grammar.xsd"
    xml:lang="en-US" version="1.0"
    root="pizza">

<rule id="digit" scope="public">
    <one-of>
        <item>one</item>
        <item>two</item>
        <item>three</item>
        <item>four</item>
        <item>five</item>
        <item>six</item>
        <item>seven</item>
        <item>eight</item>
        <item>nine</item>
    </one-of>
</rule>
<rule id="number" scope="public">
    <item repeat="1-"><ruleref uri="#digit"/></item>
</rule>
<rule id="kind">
    <one-of>
        <item>mushrum</item>
        <item>salami</item>
        <item>cheese</item>
        <item>feferony</item>
    </one-of>
</rule>
<rule id="drink" scope="public">
    <one-of>
        <item>
            <one-of>
                <item>big</item>
                <item>small</item>
            </one-of>
            colla
        </item>
        <item>juice<\item>

```

```

    </one-of>
  </rule>
  <rule id="pizza" scope="public">
    <item repeat="0-1"><ruleref uri="#number"/></item>
    <item><ruleref uri="#kind" /\></item>
    <item repeat="0-1"> pizza </item>
    and
    <item repeat="0-1"><ruleref uri="#drink"/></item>
  </rule>
</grammar>

```

### 2.4.3 Elementy SRGS

Tabulka 2.1 obsahuje přehled elementu SRGS jazyka.

Jméno	ABNF příklad	XML příklad	Popis
Pravidlo	\$jmeno =	<rule uri="#jmeno"/>	
Alternativa	A B	xml elementy na stejné úrovni	může být vysloveno A nebo B
Opakování	<1-9>	<item repeat="1-9">	určí počet možných opakování bloku
Uzavření do bloku	(vol1   vol2 )	<one-of>         <item>vol1</item>         <item>vol2</item>         </one-of>	aby celý blok mohl být na straně pravidla —
Nepovinný blok	[ vol1 ]	<one-of repeat="0-1">	
Tag	{!{ }!}	<tag> <tag>	Rozšiřuje sémantiku <sup>8</sup>
Pravděpodobnost	/0.7/	<itemrepeat-prob="0.7">	
Reference na pravidlo	\$jmeno	<ruleref uri="#jmeno">	
Vstupní pravidlo	root \$jmeno;	<grammar root="jmeno">	Identifikuje první volané pravidlo v gramatice.

Tabulka 2.1: Přehled základních elementů jazyka SRGS/ABNF

## 2.5 Hlasový dialogový systém

IVR slouží pro nahrazení člověka jako rozhodující jednotky v ústředně. Typické použití je když zavoláte na nějakou službu a z telefonního sluchátka se ozve pro objednávku zmáčknete "1", pro reklamaci "2", a tak dále. Pokud je tato interakce vykonávána automatem namísto lidmi tak firmy ušetří peníze, což je hlavní důvod jejich používání.

Pokud je zákazník nucen pro vyjádření svého přání stisknout klávesu na telefonu je použita technologie DTFM<sup>9</sup>, která převede stisk klávesy na specifický tón. Ten IVR rozezná a podle něj rozhodne o dalším postupu. Při použití ASR (Automatické rozpoznávání řeči), není uživatel nucen používat klávesnici a plynule pokračuje v rozhovoru. To přináší

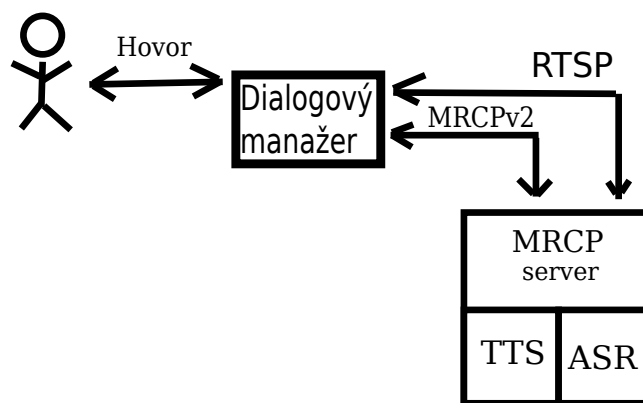
<sup>9</sup>Detaily o DTFM naleznete v RFC4733 na webu <http://tools.ietf.org/html/rfc4733>

výraznější spokojenost zákazníků (podle článku [5] až o 20%). Nemluvě o použití IVR jako vstupu pro hlasem ovládané prostředí.

Při automatickém rozhovoru může IVR přehrávat uživateli předpřipravené nahrávky nebo využít TTS (Text To Speech), převod textu na řeč.

### 2.5.1 Obecné schéma IVR systému

Je na obrázku 2.10. Uživatel vede hovor s dialogovým manažerem, který má nadefinovaný scénář, kterým uživatele provede. Po otevření telefonní linky může přehrát předpřipravenou nahrávku nebo vygenerovanou pomocí TTS. Poté pomocí protokolu MRCP inicializuje ASR rozpoznávání. Další MRCP zpráva obsahuje SRGS element, kterým informuje ASR server o tom, co očekává, že uživatel řekne, co je validní uživatelskou odpovědí. Dialogový manažer přeposílá uživatelskou odpověď na ASR server, který po nalezení shody s SRGS v těle MRCP zprávy vrací rozpoznanou promluvu v textovém řetězci. V případě vypršení časového limitu, nebo-li nerozpoznání informuje dialogový manažer chybovým stavem. Ten poté zareaguje na odpověď podle svého scénáře.



Obrázek 2.10: Schéma IVR systému.

### 2.5.2 Existující řešení

## 2.6 Protokol MRCP

Media Resource Control Protocol [3]. Je textový, síťový protokol, který je určen pro komunikaci klientů s řečovými servery. MRCP klienti jím řídí průběh třeba rozpoznávání řeči nebo syntézy řeči.

Spojení klienta a serveru se řeší pomocí SIP<sup>10</sup> (Session Initiation Protocol) hlas se přenáší pomocí RTP<sup>11</sup> (Real Time Streaming Protocol). Standardní protokoly dovolují implementaci systému do produktů třetích stran, jako je třeba FreeSWITCH<sup>12</sup> server pro IP telefonii.

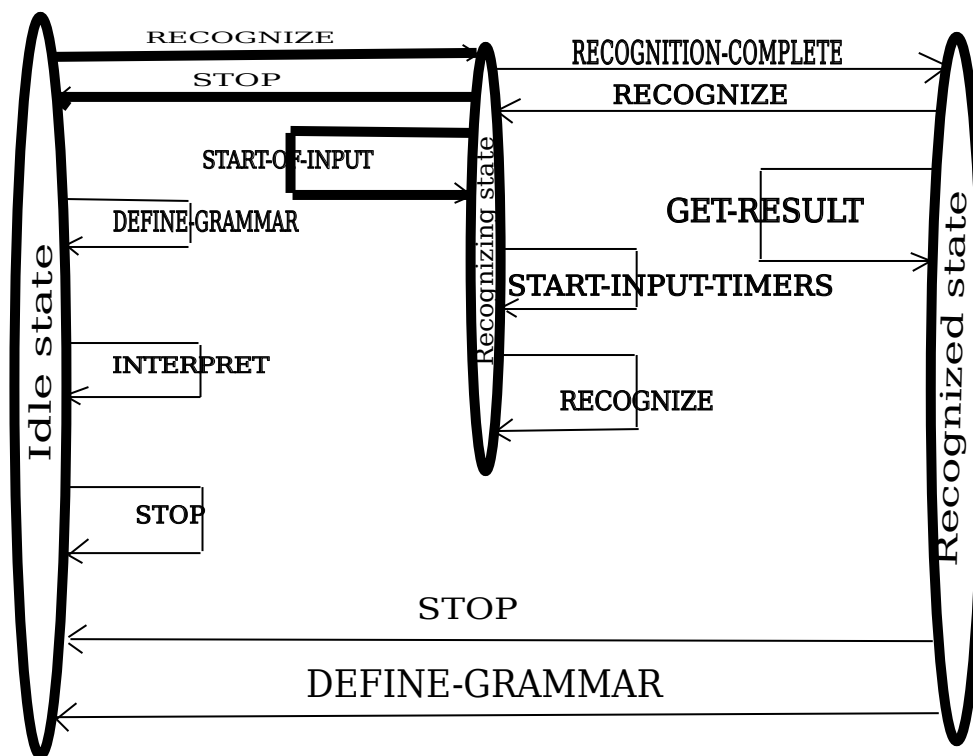
<sup>10</sup>Více v RFC3261, dostupného na webu <http://www.ietf.org/rfc/rfc3261.txt>

<sup>11</sup>Více v RFC3550, dostupného z webu <http://tools.ietf.org/html/rfc3550>

<sup>12</sup>Opensource projekt s domovským webem na adrese <http://www.freeswitch.org/>

### 2.6.1 Proces rozpoznávání v MRCP

MRCP klient řídí proces rozpoznávání na serveru podle stavového automatu zobrazeného na obrázku 2.11. Na něm je patrné že rozpoznávač nemusí po úspěšném rozpoznání vzoru, stav "Recognized", skončit.



Obrázek 2.11: Stavový automat rozpoznávání.

### 2.6.2 UniMRCP

UniMRCP<sup>13</sup> je open source projekt pro podporu vývoje MRCP serveru a MRCP klientu. Podporujícího MRCPv1 a MRCPv2. MRCPv1 používá k řízení sezení RTP a MRCPv2 SIP. Podporuje jak Linuxové tak Windowsové systémy.

Projekt je napsán v C/C++ s modulární strukturou. To dovoluje přidat novou funkcionalitu napsáním modulu, to obnáší vyplnění těla několika funkcí, které jsou definovány jako rozhraní.

UniMRCP projekt obsahuje rozhraní pro tyto MRCP aplikace:

- ASR

<sup>13</sup>Domovská stránka projektu je na <http://www.unimrcp.org/>

- TTS
- SR (Speech Recorder), nahrávání řeči.
- SVI (Speaker Verifier and Identifier), ověření a identifikace řečníka.

Projekt se především zaměřuje na podporu implementace MRCP klienta a nesnaží se vytvořit řečový server. UniMRCP obsahuje ukázkové implementace rozhraní pro MRCP aplikace.



## Kapitola 3

# Převod SGRS gramatik na váhovaný konečný automat

V této kapitole bude vysvětleno jak je SGRS/ABNF gramatika převedena na konečný automat reprezentovaný pomocí tříd z knihovny OpenFST. Ke čtení zapsané gramatiky je využita knihovna ANTLR<sup>1</sup> (ANother Tool for Language Recognition), která generuje programový kód v C/C++ pro lexikální a syntaktickou analýzu specifikovaného jazyka.

Pro implementaci převodu byla zvolena sada elementů SGRS/ABNF, která dovoluje základní práci s gramatikami. Ostatní prvky jazyka podle [2] jsou rozpoznány, ale nejsou použity. Uživateli je ovšem o tomto faktu předaná zpráva.

Seznam elementů SGRS/ABNF, které se projeví na tvaru výstupního váhovaného konečného automatu (VKA):

### Slovo

Jednotlivá rozpoznávaná slova tvoří terminály v SGRS. Obrázek 3.1. V implementaci musí být zapsána velkými písmeny.



Obrázek 3.1: Schéma elementu slovo v VKA G.

### Definice pravidla

S předpisem `$rulename=R;` se uloží pod jménem `rulename` do seznamu pravidel.



Obrázek 3.2: Schéma elementu pro definování pravidla v VKA G.

---

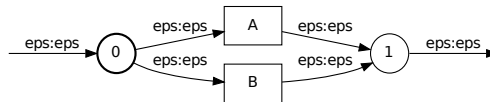
<sup>1</sup>Domovská stránka projektu je webu <http://www.antlr.org/>

### Reference na pravidlo

s předpisem  $\$R$ . Ze seznamu pravidel zkopíruje jeho graf a vloží ho na požadované místo definici nového pravidla.

### Alternativa

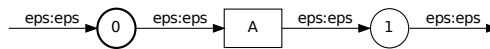
S předpisem  $A|B$  se do pravidla gramatiky vloží, tak že se vedou  $\varepsilon$  přechody od obou možností průchodů. Detail na obrázku 3.3. A a B představují podgrafy vytvořené ze zápisu gramatiky pro volbu A nebo B.



Obrázek 3.3: Schéma elementu pro výběr alternativy v VKA G.

### Blok

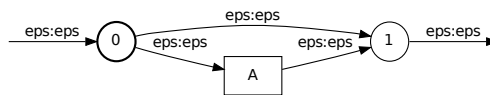
Uzavírá část A do bloku obaleného stavu, aby na něj bylo možno aplikovat elementy SRGS. Obrázek 3.4.



Obrázek 3.4: Schéma elementu pro vybrání jednu z možností nebo nic v VKA G.

### Nepovinný blok

Blok rozšiřuje rozšiřuje o  $\varepsilon$  hranu, která ho dovoluje přeskočit. A tak vytvořit nepovinnou část gramatiky. Obrázek 3.5.

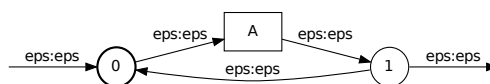


Obrázek 3.5: Schéma elementu pro vybrání jednu z možností nebo nic v VKA G.

### Opakování

S předpisem  $\langle 1 \rightarrow \rangle$  znamená, že část gramatiky se může do nekonečna opakovat. Realizuje zpětnou  $\varepsilon$  hranou. Obrázek 3.6. Jednička může být nahrazena jakýmkoliv kladným číslem.

Konkrétní počet opakování s předpisem  $\langle 4 \rangle$  se realizuje tak, že se opakovaný blok nakopíruje za svůj konec.

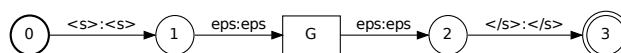


Obrázek 3.6: Schéma elementu pro nekonečné opakování v VKA G.

Rozsah s předpisem <2-4> je pak realizován kombinací konkrétního počtu opakování a nepovinných bloků.

### Vstupní pravidlo

S předpisem `root=$A;`, vyhledá vstupní pravidlo do gramatiky, které rozšíří aby tvořilo plnohodnotný automat. Obrázek 3.7.



Obrázek 3.7: Schéma elementu pro začátek gramatiky v VKA G.

## 3.1 Vstup

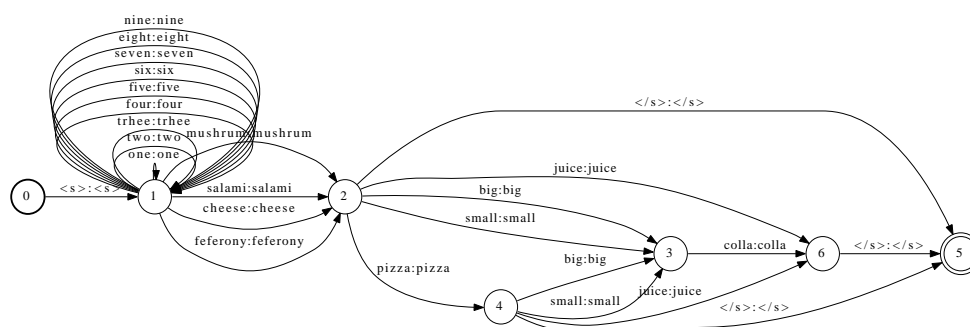
Vstupem do převodu je SRGS/ABNF gramatika. Příklad gramatiky zpracování požadavku v pizzerii: `#ABNF 1.0 ISO-8859-1;`  
`language en-US;`  
`// root role is start point of gramar`  
`root $pizza;`  
`public $digit = one | two | three | four | five | six | seven | eight | nine;`  
`public $number= $digit <1-> ;`  
`public $kind = mushrum | salami | cheese | feferony;`  
`public $drink = (big | small) colla | juice;`  
`public $pizza = [ $number ] $kind [ pizza ] [ $drink ];`

## 3.2 Výstup

Výstupem z převodu je váhovaný konečný automat. Obrázek 3.8 je výstupem k příkladu v předešlé kapitole 3.1.

## 3.3 Postup převodu gramatiky na graf G

K analýze SRGS/ABNF gramatiky je použit nástroj ANTLR podle knihy [10]. Po identifikaci každého z přijímaných elementu z SRGS/ABNF se tento element převede na část grafu v OpenFST. Pokud je identifikován element, který převodník ignoruje, tak se jde dál. Analyzátor tedy dokáže zpracovat celý standart SRGS/ABNF aniž by skončil chybou.



Obrázek 3.8: Vizualizace převodníku G (zpracované gramatiky) v OpenFST.

Po převedení gramatiky na konečný automat se k automatu přidá stav s označením 0, který má orientovanou hranu vedoucí z prvního stavu v automatu. Tato hrana má na svém vstupu i výstupu symbol `<s>`. Stejně tak je přidán koncový stav s orientovanou hranou vedoucí do posledního stavu automatu. Tato nová hrana má za svůj vstupní i výstupní symbol `</s>`. Tyto symboly označují začátek a konec rozpoznané sekvence slov.

Po vytvoření VKA jsou aplikovány operace determinizace a minimalizace, aby se zmenšila jeho velikost a vznikl VKA bez  $\varepsilon$  přechodů.

### 3.4 ANTRL

Je nástroj napsaný v Javě pro generování kódu syntaktické a lexikální analýzy nadefinovaného doménového<sup>2</sup> jazyka. Kódy mohou být kromě Javy generovány do jazyků jako Python a C. Konkrétně generátoru pro jazyk C bylo využito v tomto projektu.

Výhodou ANTRL je menší množství kódu které musí programátor pro zpracování vstupu vytvořit. A důsledná kontrola vstupu, zda odpovídá nadefinovanému schématu.

<sup>2</sup>Takto je označován jazyk pro specifické použití.

Další vlastností, kterou ale implementace pro BSAPI nepoužívá, je tvorba abstraktního syntaktického stromu.

### 3.4.1 ANTRL pravidla pro SGRS/ABNF

Se dělí na lexikální a sémantická pravidla. Jména lexikálních pravidel začínají velkými písmeny a jména syntaktických pravidel zase malými. Za jménem následuje znak `:` za ním je definice pravidla. Řetězce v apostrofech jsou definicí lexému.

Lexikální analýza očekává vstup ve formátu utf-8. SRGS určuje že se pro definici číslic a znaků mají použít stejné symboly jako v XML<sup>3</sup>. Ukázka několika dalších definovaných lexikálních pravidel:

```
WS : ( ' ' | '\r' | '\t' | '\u000C' | '\n' ) { $channel=HIDDEN; } ;

fragment Letter : BaseChar | Ideographic;

fragment NameChar:
    Letter | Digit | '.' | '-' | '_' | ':' | CombiningChar |
    Extender
;
fragment Name:
    (Letter | '_' | ':') (NameChar)*
;
fragment Names:
    Name (WS Name)*
;
fragment Nmtoken:
    (NameChar)+
;

fragment Nmtokens:
    Nmtoken (WS Nmtoken)*
;

DoubleQuotedCharacters :
    '"' ( '\u0000' .. '\u0021' | '\u0023' .. '\uFFFE' ) * '"'
;

Token:
    (Letter)+
    | DoubleQuotedCharacters
;

RuleName:
    '$' Name
;

ML_COMMENT : '/*' ( options {greedy=false;} : . ) * '*/' { $channel=
    HIDDEN; };
SL_COMMENT : '//' ( options {greedy=false;} : . ) * '\n' { $channel=
    HIDDEN; };
```

---

<sup>3</sup>Seznam povolených symbolů na webu <http://www.w3.org/TR/2000/REC-xml-20001006#NT-Letter>

Klíčové slovo **fragment** jazyka ANTRL říká, že definované lexikální pravidlo netvoří lexémy, ale dá se pouze použít jako část dalšího lexikálního pravidla.

Sémantická pravidla pak kombinují lexikální pravidla, lexémy a jiná sémantická pravidla. Základní sémantická pravidla pro jazyk SRGS/ABNF v jazyku ANTRL:

```
mainRule:
  declaration* (( ruleDefinition | sequence ) ) * EOF
;

declaration : rootRuleDecl;
rootRuleDecl: 'root' RuleName ';' ;

ruleDefinition: ruleScope? ruleName=RuleName '=' ruleExpansion ';' ;

ruleScope: 'private' | 'public';

ruleExpansion : ruleAlternative ( '|' ruleAlternative ) * ;

sequence: (sequenceElement)+ ';' ;

ruleAlternative: (sequenceElement)+ ;

blockExpansion : blockAlternative ( '|' blockAlternative ) * ;

blockAlternative: (sequenceElement)+ ;

sequenceElement: subexpansion (repeatOperator)? ;

subexpansion:
  readToken
  | ruleRef
  | tag
  | '(' choiceOne ')' languageAttachment?
  | '[' optionalOne ']' languageAttachment?
  /* 0 or 1 */
;

readToken:
  Token
  languageAttachment?;

choiceOne: blockExpansion;

optionalOne: blockExpansion;

ruleRef: localRuleRef | specialRuleRef | externalRuleRef;

localRuleRef: RuleName ;

specialRuleRef: '$NULL' | '$VOID' | '$GARBAGE';

externalRuleRef: '$' '<' Token '>' ('~' '<' Token '>')? ;
```

```

tag: tagPos=('{' .* '}' | '{!{' .* '}!}' ) ;

repeatOperator:
  '<' Number ( '-' )? ( Number )? '>' ;

languageAttachment: '!' .* ';' ;

```

Vstupem pro sémantickou analýzu gramatiky je pravidlo `mainRule`, které v implementaci vrací seznam vygenerovaných stavů a hran mezi nimi. Podle vrácených dat se sestaví síť převodníku G za pomoci knihovny `OpenFst`.

### 3.4.2 Zpracování chyb

Pokud je při analýze objevena chyba ve vstupu je zavolána předdefinovaná metoda pro její výpis uživateli. Na standardní chybový výstup vytiskne dosti nesrozumitelné hlášení. Toto chování se dá obejít náhradou této funkce.

Program pro analýzu SRGS/ABNF gramatik chybový callback přepisuje funkcí `myDisplayRecognitionError`. Vstupem je struktura<sup>4</sup> `ANTLR3_BASE_RECOGNIZER`, která nese informace o současném stavu průběhu analýzy. Funkce chybu zpracuje a uloží, aby mohla být předána do BSAPI. Její implementace je zde:

```

void myDisplayRecognitionError
(pANTLR3_BASE_RECOGNIZER grammarAnalyser,
 pANTLR3_UINT8 * tokenNames)
{
    //grammarAnalyser->state->exception->token
    long line      = grammarAnalyser->state->exception->line;
    long position = grammarAnalyser->state->exception->
        charPositionInLine;

    char * errorText = (char*)grammarAnalyser->state->exception->name;

    switch(recognizer->state->exception->type)
    {
        case ANTLR3_RECOGNITION_EXCEPTION:
            break;
        case ANTLR3_MISMATCHED_TOKEN_EXCEPTION:
            break;
        case ANTLR3_NO_VIABLE_ALT_EXCEPTION:
            break;
        case ANTLR3_MISMATCHED_SET_EXCEPTION:
            break;
        case ANTLR3_EARLY_EXIT_EXCEPTION:
            break;
        case ANTLR3_FAILED_PREDICATE_EXCEPTION:
            break;
        default:
            break;
    }
}

```

---

<sup>4</sup>Dokumentace k struktuře reprezentující chybu během analýzy je na webu [http://antlr3.org/api/C/struct\\_a\\_n\\_t\\_l\\_r3\\_\\_\\_e\\_x\\_c\\_e\\_p\\_t\\_i\\_o\\_n\\_d\\_struct.html](http://antlr3.org/api/C/struct_a_n_t_l_r3___e_x_c_e_p_t_i_o_n_d_struct.html)

V definici funkce je vidět rozdělení chyb na několik typů. Kde třeba typ `ANTLR3_MISMATCHED_TOKEN_EXCEPTION` znamená, že další načtený lexém neodpovídá žádné možnosti, která je definována pravidlem.

Uživatelsky nejpříjemnější metodou zpracování chyb by bylo potlačení callbacku (odkaz na funkci) a analyzování abstraktního syntaktického stromu, který usnadňuje zotavení z chybového stavu a pokračování v analýze.



## Kapitola 4

# Sestavení rozpoznávací sítě

Jádrem sestavení rozpoznávací sítě je vzorec  $HCLG = H \circ \min(C \circ \det(L \circ G))$  vyplývající z článku [7]. Pro praktickou aplikaci s knihovnou OpenFst, je potřeba vzorec rozšířit na:

$$LG = \text{arcsort}_{\text{olabel}}(\text{rmepsilon}(L)) \circ \text{arcsort}_{\text{ilabel}}(\text{rmepsilon}(G)) \quad (4.1)$$

$$CLG = \text{arcsort}_{\text{olabel}}(\text{rmepsilon}(C)) \circ \text{arcsort}_{\text{ilabel}}(\det(\text{rmepsilon}(LG))) \quad (4.2)$$

$$HCLG = \text{arcsort}_{\text{olabel}}(H) \circ \text{arcsort}_{\text{ilabel}}(\min(\det(\text{rmepsilon}(CLG)))) \quad (4.3)$$

Tento vzorec pak vytvoří výstup v podobě převodníku, kde každá hrana má svůj vstupní symbol, výstupní symbol, a pravděpodobnost. Situaci zachycuje výřez  $HCLG$  převodníku vytvořeného v OpenFst:

```
...
3    4    sil_4    eps 0.0651507005
4    2    sil_2    eps 0.105728
4    4    sil_4    eps 0.00762159005
4    5    eps      eps 0.0566494986
5    6    z_s2_2    ZERO 2.30273438
5    7    s_s2_1    SEVEN 2.30273438
...
```

První dva sloupce jsou ID patřící uzlům v grafu. První je uzel, ze kterého vede hrana, a druhý je uzel, do kterého vede hrana. Třetí sloupec je akustická jednotka čtvrtý sloupec je slovo. A jako poslední je uložena pravděpodobnost přechodu přes hranu.

### 4.1 Format rozpoznávací sítě v BSAPI

Který dekodér načítá je unikátní pro projekt. Od převodníku v OpenFst se liší tím, že má symbol z hrany převeden do uzlu. Proto se při sestavení rozpoznávací sítě v BSAPI provádí několik kroků navíc:

- Převod na akceptor  $HCLG$ a se slovními a akustickými uzly.
- Převod na formát pro BSAPI.
- Seřazení grafu podle počátečního uzlu a přechod na numerické číslování uzlu.

Tyto kroky umožňují snadný zásah a úpravu výsledné rozpoznávací sítě. Pro případ doplnění speciální funkcionality dekodéru, která mu může být předána skrze síť.

BSAPI podporuje binární i textový formát sítě. Textový dovoluje snadnější vývoj a kontrolu zásahu do sítě. Je podporován i z historických důvodů, protože ho BSAPI používalo před binárním a proto že dosavadní metody, tvoření sítí, byly postaveny nad linuxovými skripty. Výhodou binárního formátu sítě je zase rychlost, se kterou ji dekodér načítás.

#### 4.1.1 Převod na akceptor *HCLGa*

VKA je převodník s totožným vstupním i výstupním symbolem. Aby dekodér, který bude používat rozpoznávací síť rozeznal, který symbol je akustický model a který je slovo jsou symboly rozšířené o předponu "W=" nebo "M=". "W=" jako "word" je předpona pro symboly z G a "M=" jako model je předpona pro symboly z H. Například:

```
0      1      W=<s>    W=<s>
1      2      M=sil_2 M=sil_2 -2.30258512
2      2      M=sil_2 M=sil_2 0.00865190011
...
```

Druhý a třetí sloupec pro vstupní a výstupní symboly obsahuje stejné hodnoty. Toho lze docílit algoritmem 3. Kde:

- . lze představit jako operaci konkatenace dvou řetězců.
- && je symbolem pro logické and, stejně jako v podmínce jazyka C++.

**Data:** Převodník  $HCLG = (A, B, Q_T, I, F_T, E, \lambda, p_T)$  s jedním počátečním stavem.

**Result:** Akceptor  $HCLGa = (Q_A, \Sigma, \sigma, s, F_A, p_A)$

method

$s \leftarrow I$

$F_A \leftarrow F_T$

$p_A \leftarrow p_T$

$Q_A \leftarrow Q_T$

**while**  $Q_T \neq \emptyset$  **do**

$state \leftarrow HEAD(Q_T)$

$DEQUEUE(Q_T)$

**for each**  $h, h \in E \ \&\& \ h \in (state, input_l, output_l, likelihood, next_s)$  **do**

$\Sigma \leftarrow \Sigma \cup "M = " \cdot input_l, "W = " \cdot output_l$  **if**  $input_l = \varepsilon \ \&\& \ output_l = \varepsilon$

**then**

$\sigma \leftarrow \sigma \cup (state, \varepsilon, likelihood, next_s)$

**end**

**if**  $input_l = \varepsilon \ \&\& \ output_l \neq \varepsilon$  **then**

$\sigma \leftarrow \sigma \cup (state, "W = " \cdot output_l, likelihood, next_s)$

**end**

**if**  $input_l \neq \varepsilon \ \&\& \ output_l = \varepsilon$  **then**

$\sigma \leftarrow \sigma \cup (state, "M = " \cdot input_l, likelihood, next_s)$

**end**

**if**  $input_l \neq \varepsilon \ \&\& \ output_l \neq \varepsilon$  **then**

$Q_A \leftarrow \cup new_s$

$\sigma \leftarrow \sigma \cup (state, "M = " \cdot input_l, likelihood, new_s)$

$\sigma \leftarrow \sigma \cup (new_s, "W = " \cdot output_l, 0, next_s)$

**end**

**end**

**end**

**Algorithm 3:** Metoda převodu převodníku T na akceptor A.

Dalšími kroky v převodu je determinizace podle [7], kterou lze na  $HCLGa$  provést díky algoritmu 1. Pak následuje přemapování pomocných rozlišovacích symbolů na  $\varepsilon$ . Pomocné symboly začínající znakem #, které slouží k rozlišení stejně znějících slov v převodníku L, tak aby šel determinizovat. Tedy můžeme zapsat:

$determinize(map(\#?, \varepsilon))$

#### 4.1.2 Převod na formát pro BSAPI

Převodem na mezisítě se usnadní převod na rozpoznávací síť. Formát mezi sítě je zachycen zde:

```
n0 0 l=-0
0 W=<s> n1
n1 1 l=-0
1 M=sil_2 n2
n2 2 l=-0.0579227
2 M=sil_4 n4
n2 3 l=-0.0861414
3 M=sil_3 n3
n2 1 l=-0.0086519
n3 2 l=-0.0651507
...
```

```
n422 n425 l=-0.0566495
...
```

Jak je vidět formát BSAPI obsahuje 3 typy záznamu. První je spojení mezi uzly `n1 1 l=-0`, druhý je spojení se symbolem `0 W=<s> n1` a třetí je spojení ke koncovému uzlu `n422 n425 l=-0.0566495`. Lépe to vysvětluje následující výřez z poznámek ve zdrojových kódech pro uchování tohoto formátu:

```
* Which describe node with label and connetion
* to other nodes with path likelihood.
* One record (describe or connection) is on one line.
* node with label syntax example:
* "0 W=<s> n1"
* | | '- mToNode = 1, is pointer to node1
* | '----- mNodeLabel = "W=<s>"
* '----- mFromNode (node id ) = 0
* mLikelihood = 0.0 (implicit)
*
* and connection to other node
* "n2 2 l=-0.0687369"
* | | '----- likelihood = -0.0687369
* | '----- toNode = 2
* '----- fromNode = 2
* nodeLabel = "" (implicit)
*
* and connection to other node (end node)
* "n194 n195 l=-0.0614271"
* | | '----- likelihood = -0.0614271
* | '----- toNode = 195
* '----- fromNode = 194
* nodeLabel = "" (implicit)
```

Takovéhoto formátu lze docílit algoritmem 4, kde funkce *NMAP* mapuje string na unikátní ID a pro stejný string vrátí stejné ID. Pro jeho potřeby si mezisít TNET zadefinujeme jako trojici nad okruhem  $K$  s reálnými čísly, který je použit pro vyjádření pravděpodobnosti:

$TNET = (Q, \sum, \delta)$ , kde:

$Q$  je konečná množina stavů.

$\sum$  je konečná množina symbolů.

$\delta$  je množina přechodů,

taková že  $\delta \subseteq Q \times \sum \times K \times type \in \{Label, Connection, ConnectionToEnd\} \times Q$  kde první stav značí startovní symbol přechodu a druhý stav značí koncový stav přechodu.

**Data:** Deterministický akceptor  $HCLGa = (Q_A, \sum_A, \sigma, s, F, p)$ , bez  $\varepsilon$  symbolů.

**Result:**  $TNET = (Q, \sum, \delta)$

**Method:**

$\sum \leftarrow \sum_A$

$\delta \leftarrow \{\varepsilon\}$

$Q \leftarrow Q_A$

**forall the**  $(from_s, symbol, likelihood, to_s) \in \sigma$  **do**

**if**  $symbol = \varepsilon$  **then**

$\delta \leftarrow \delta \cup \{(from_s, \varepsilon, 0 - likelihood, Connection, to_s)\}$

**end**

**if**  $symbol \neq \varepsilon$  **then**

$m \leftarrow symbol \cdot \_ \cdot ToString(to_s)$

$Q \leftarrow Q \cup NMAP\{(m)\}$

$\delta \leftarrow \delta \cup \{(from_s, \varepsilon, 0 - likelihood, ConnectionToEnd, NMAP(m))\}$

**if**  $to_s \notin S$  **then**

$S \leftarrow S \cup \{to_s\}$

$\delta \leftarrow \delta \cup \{(NMAP(m), symbol, 0, Label, to_s)\}$

**end**

**end**

**end**

**Algorithm 4:** Převod na formát pro BSAPI.

#### 4.1.3 Optimalizace BSAPI sítě

Obnáší setřídění podle počátečního uzlu a numerické setřídění uzlů.

Formát textového souboru s optimalizovanou BSAPI sítí:

```
N=467 L=946
0 W=!NULL 1 l=0
1 W=<s> 2 l=0
2 M=sil_2 3 l=-0.05792 4 l=-0.08614 2 l=-0.00865
3 M=sil_4 5 l=-2.35938 6 l=-2.35938 7 l=-1.66623
8 l=-2.35938 9 l=-2.35938 10 l=-2.35938 11 l=-1.66623
12 l=-2.35938 3 l=-0.00762 2 l=-0.10572
...
```

N na prvním řádků udává počet uzlů v síti a L počet hran v síti.

Pro potřeby algoritmu 5, který setřídí síť zadefinujeme objekt  $NET$  nad okruhem  $K$ .

Okruh slouží pro vyjádření pravděpodobností přechodů ze stavu do stavu:

$NET = (Q, \sum, \varrho, \phi)$ , kde:

$Q$  je konečná množina stavů sítě.

$\sum$  je konečná množina symbolů.

$\varrho$  je funkce mapující symbol na stav  $\varrho : Q \rightarrow \sum$

$\phi$  je množina přechodů ze stavu do stavu, taková že  $\phi \subseteq Q \times K \times Q$

**Data:**  $TNET = (Q_T, \sum_T, \delta)$   
**Result:**  $NET = (Q, \sum, \varrho, \phi)$   
**Method:**  
**while**  $Q_T \neq \emptyset$  **do**  
     $n \leftarrow HEAD(Q_T)$   
     $DEQUEUE(Q_T)$   
     $Q \leftarrow Q \cup \{n\}$   
    **for** *each*  $(from_s, symbol, likelihood, ConnectionType, to_s) \in \delta$  **do**  
        **if**  $from_s = n$  **and**  $ConnectionType = Label$  **then**  
             $\varrho \leftarrow \{(x, -) \in \varrho\}$ , where  $x \neq n$   
             $\varrho \leftarrow \varrho \cup \{(n, symbol)\}$   
        **end**  
        **if**  $from_s = n$  **and**  $ConnectionType \neq Label$  **then**  
             $\phi \leftarrow \phi \cup \{(n, likelihood, to_s)\}$   
        **end**  
    **end**  
**end**

**Algorithm 5:** Optimalizace BSAPI sítě.

#### 4.1.4 Vytvoření převodníku L

Převodník L se tvoří každé gramatice na míru. Protože převodník obsahující všechna slova v jazyku by byl příliš velký a s takovým převodníkem by se během sestavování rozpoznávací sítě pracovalo pomalu. Navíc všechna slova v daném jazyku ani obsahovat nemůže.

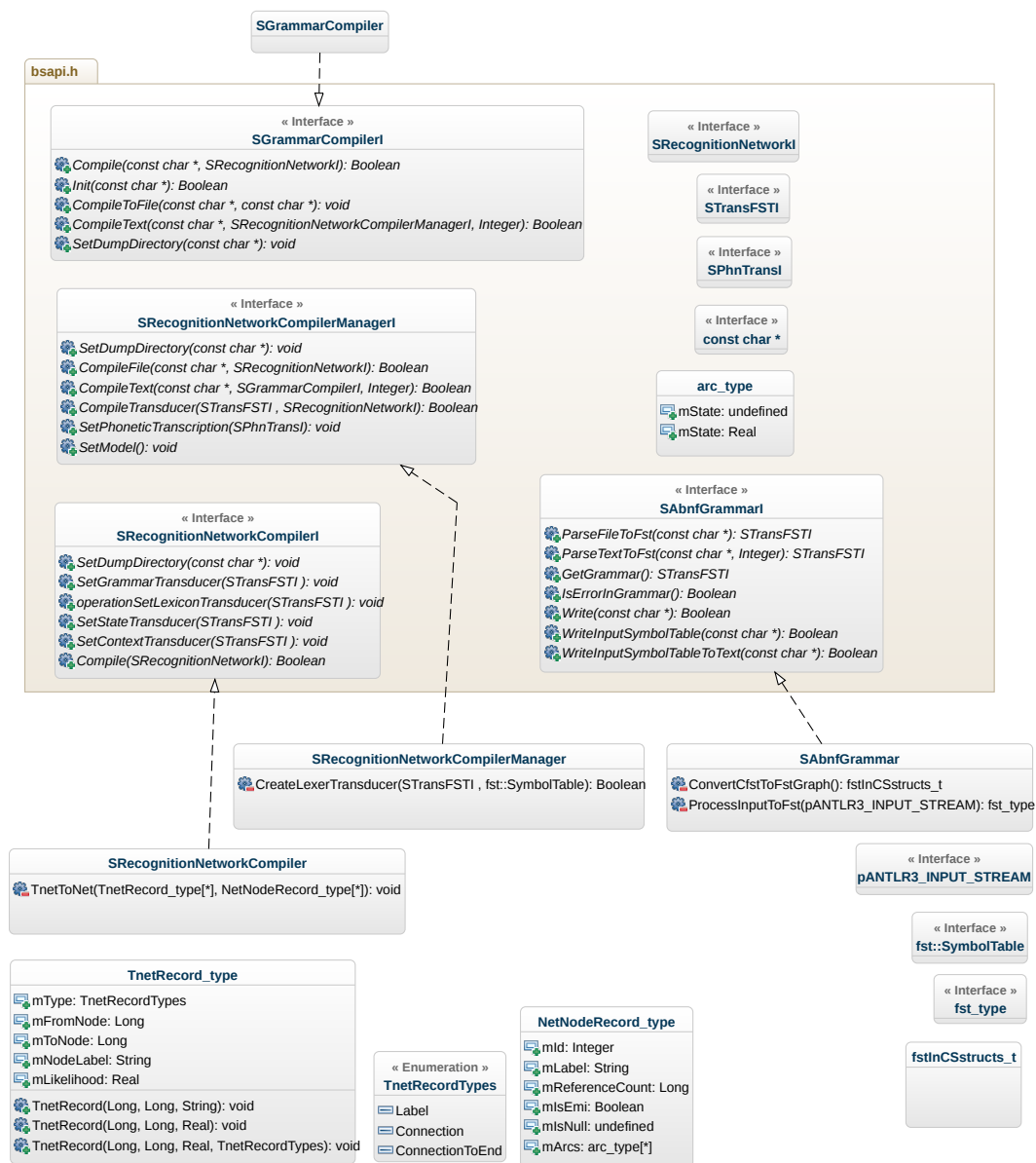
V BSAPI slouží pro překlad slov na fonémy rozhraní SPhnTrasI, které vrací fonémový přepis slova. To se nejdříve pokusí vyhledat slovo ve slovníkům známých výslovností a pokud ho nenajde tak vrátí přepis grafému na fonémy. Grafém je znak z abecedy. Přepis je definován pro každý jazyk. Přepis grafémů vrací nesprávnou výslovnost především pro slova, která jazyk přejal z jiného jazyku a jenž vyslovuje variantou originální výslovnosti.

## 4.2 Zapouzdření do BSAPI

Obrázek 4.2 je UML (Unified Modeling Language) diagramem tříd zachycující objekty zpracovávající SRGS/ABNF gramatiku a objekty pro sestavení rozpoznávací sítě. V diagramu je znázorněn soubor bsapi.h, ve kterém jsou definovány všechny rozhraní, které knihovna používá pro komunikaci. Také to je hlavičkový soubor, který se dodává zákazníkům. Jednotlivé objekty blíže vysvětluje tabulka 4.2.

Jméno objektu	Popis objektu
SGramarCompilerI	rozhraní kompilátoru sítě.
SGramarCompiler	Implementace rozhraní SGramarCompilerI.
SRecognitionNetworkI	Rozhraní pro práci s rozpoznávací sítí.
SRecognitionNetwork	Implementace rozhraní SRecognitionNetworkI.
STransFSTI	Rozhraní pro práci s konečnými automaty.
SPhnTrasI	Rozhraní pro přepis řetězce reprezentující slovo na jeho fonémy. Používá se při sestavování převodníku L.
arc_type	Struktura popisující hranu sítě, Obsahující cílový stav a pravděpodobnost hrany.
SRecognitionNetworkCompilerManagerI	Rozhraní pro přípravu převodníků pro sestavení rozpoznávací sítě.
SRecognitionNetworkCompilerManager	Implementace rozhraní SRecognitionNetworkCompilerManagerI.
SAbnfGrammarI	Rozhraní pro načítání a zpracování SRGS/ABNF gramatiky
SAbnfGrammar	Implementace rozhraní SAbnfGrammarI.
SRecognitionNetworkCompilerI	Rozhraní pro sestavení rozpoznávací sítě dle kapitoly 4
SRecognitionNetworkCompiler	Implementace rozhraní SRecognitionNetworkCompilerI

Tabulka 4.1: Zodpovědnosti jednotlivých objektů.



Obrázek 4.1: Zapozdření tříd pro sestavení rozpoznávací sítě do BSAPI.



## Kapitola 5

# Součinnost gramatik s rozpoznávačem

Funguje skrze rozpoznávací sítě, které jsou vytvořeny podle gramatiky a poté načteny dekodérem rozpoznávače. Rozpoznávač poté hledá nejlepší průchod rozpoznávací sítí, tak aby odpovídal zpracovávané řeči.

### 5.1 Testy sestavení rozpoznávací sítě

Jsem prováděl, protože mnou sestavené rozpoznávací sítě se liší od původní sítě vytvořené skriptem. Při vytváření sítí jsem totiž neimplementoval poslední krok, kterým je optimalizace sítě nástrojem `SExpand` z projektu <sup>1</sup>.

Testy jsem prováděl na stroji s procesorem i5-3470<sup>2</sup>.

Přesnost, rychlost a paměťové nároky rozpoznávací sítě jsem testoval na sítích pro přepis řeči, kde gramatika dovoluje všechny kombinace slov. Více v o statistickém popisu jazykových modelů v článku [1] a [8]. Tato úprava mi dovolí využít připravené testy nad sadami nahrávek pro testování jazykových modelů. Díky tomu jsem dostal objektivnější výsledky, protože jsem je mohl srovnat s předešlou metodou tvoření rozpoznávací sítě. Testy byly provedeny na testovacích sadách pro CZ (český jazyk) a RU (ruský jazyk).

#### 5.1.1 Rychlost načtení gramatiky a sestavení rozpoznávací sítě

Byla testována na operačním systému Linux s pomocí knihovny `<sys/time.h>` z prostředí jazyka C++. Před a za kód pro zpracování vstupu a získání rozpoznávací sítě jsem vložil kód pro zaznamenání času.

Doba zpracování vstupu je závislá na jeho velikosti, počtu pravidel a na jejich vzájemném propojení.

Ukázková gramatika pro objednání pizzy 3.1 se zpracovávala do rozpoznávací sítě zhruba za 0,158 s. Tento čas je výrazně kratší než 3 minuty, které byly zapotřebí pro sestavení souboru obsahující ekvivalentní rozpoznávací síť starou metodou vytvořenou ve skriptech.

Do měřeného času není započítán čas pro inicializaci objektu pro `SRecognitionNetworkCompilerManager`, která se provádí jen jednou. Po inicializaci se

<sup>1</sup>více o STK na webu <http://speech.fit.vutbr.cz/cs/software/hmm-toolkit-stk>

<sup>2</sup>Procesor Intel(R) Core(TM) i5-3470 CPU @ 3.20GHz, více na webu [http://ark.intel.com/products/68316/Intel-Core-i5-3470-Processor-6M-Cache-up-to-3\\_60-GHz](http://ark.intel.com/products/68316/Intel-Core-i5-3470-Processor-6M-Cache-up-to-3_60-GHz)

Síť	WER [%]	WA [%]	WC [%]
CZ ze skriptu	35,2	67,2	64,8
CZ implementace	35,2	67,2	64,8
RU ze skriptu	52,6	51,5	47,4
RU implementace	52,6	51,5	47,4

Tabulka 5.1: Přesností zpracování.

může sestavení rozpoznávací sítě z gramatiky libovolně krát opakovat. U staré metody bylo zapotřebí po každém spuštění zapotřebí znovu načíst vstupní modely.

### 5.1.2 Přesnost rozpoznávací sítě

Provádím, abych zjistil jak vynechání optimalizace rozpoznávací sítě nástrojem **SExpand** ovlivnilo přesnost rozpoznávání s touto sítí. Pro vyhodnocení přesnosti rozpoznávání řeči se používají testovací sady nahrávek s referenčním přepis. Ve výsledcích testovaného rozpoznávače řeči, výsledkem je její přepis, se pro každou nahrávku sledují tyto parametry:

R (References) počet referenčních slov.

C (Correct) počet správně rozpoznávaných slov na pozici referenčních slov.

S (Substitution) počet špatně rozpoznávaných slov.

I (Insertions) počet vložených slov, mimo místa s referenčními slovy.

D (Deletions) počet míst, kde je referenční slovo, ale rozpoznávač nic nerozpoznal.

Podle článků [9] se hlídají ukazatele jako WER (Word Error Rate), WA (Word Accuracy), WC (Word Correctness). Které se vypočítají následovně:

$$WC = \frac{C}{R} \times 100 \quad (5.1)$$

$$WER = \frac{S + D + I}{R} \times 100 \quad (5.2)$$

$$WA = 1 - WER \quad (5.3)$$

Z tabulky 5.1 je vidět, že implementace skládání rozpoznávacích sítí je generuje ekvivalentní výstupy k předešlému vytváření sítí pomocí skriptů.

### 5.1.3 Rychlost rozpoznávací sítě

Provádím, abych zjistil jak vynechání optimalizace rozpoznávací sítě nástrojem **SExpand** ovlivnilo rychlost rozpoznávání s touto sítí.

Z tabulky 5.2 pro češtinu je poměr součtu dob zpracování nahrávek k součtu dob trvání nahrávek pak dává u původní metody tvorby rozpoznávacích sítí hodnotu 0,41 a u současné implementace v BSAPI pak 0,43.

Z tabulky 5.3 pro ruštinu je poměr součtu dob zpracování nahrávek k součtu dob trvání nahrávek pak dává u původní metody tvorby rozpoznávacích sítí hodnotu 0,48 a u současné implementace v BSAPI pak 0,49.

Nahrávka	Délka [s]	Původní [s]	Nový [s]
wav_prom060905_0003.wav	20,92	8,87	9,15
wav_prom060907_0006.wav	80,39	35,67	37,64
wav_prom060907_029a.wav	31,16	13,17	13,39
wav_prom060907_032a.wav	42,95	17,24	17,63
wav_prom060907_111a.wav	75,21	29,21	30,51

Tabulka 5.2: Časy zpracování nahrávek pro CZ jazykový model.

Nahrávka	Délka [s]	Původní [s]	Nový [s]
wav_ru3467.wav	74,34	36,71	37,60
wav_ru6535.wav	72,92	28,34	28,94
wav_ru6723b.wav	136,00	61,73	63,19
wav_ru7028.wav	61,91	32,90	33,84
wav_ru7914.wav	65,52	23,34	23,98
wav_ruN0001.wav	30,40	18,22	18,52
wav_ruN0002.wav	51,77	35,09	35,87

Tabulka 5.3: Časy zpracování nahrávek pro RU jazykový model.

Tato metoda měření není přesná, protože nezohledňuje množství řeči v nahrávce, která je zpracovávána. Ale přesto nabízí srovnání obou postupů a ukazuje, že současná implementace je pomalejší. Nárůst času potřebný pro zpracování nahrávky je sice malý, ale nezanedbatelný, protože se projeví při práci s velkým množstvím dat.

#### 5.1.4 Paměťové nároky rozpoznávací sítě

Provádím, abych zjistil jak vynechání optimalizace rozpoznávací sítě nástrojem **SExpand** ovlivnilo množství paměti potřebné pro rozpoznávací síť. Jsem sledoval linuxovým nástrojem **top**, který zobrazuje množství operační paměti, kterou konzumuje proces. Výsledky jsou zachyceny v tabulce 5.4, kde je zaznamenána velikost souborů pro ekvivalentní rozpoznávací síť. Jedna vyrobena původním skriptem a druhá vytvořená kompilátorem.

Dále je v tabulce 5.4 zaznamenána největší vytíženost operační paměti RAM (Random Access Memory), kterou zabírá proces rozpoznávače. Velikost obsazené paměti rozpoznávačem je ovlivněna především velikostí sítě a velikostí nahrávky. Proto se tento údaj musí chápat jako srovnávací, protože jediným rozdílem v podmínkách testů byla vstupní síť s jazykovým modelem pro rozpoznávání.

Z tabulky 5.4 je vidět že nová implementace sestavení rozpoznávacích sítí generuje větší paměťové struktury.

Jazyk	Velikost původní sítě	Velikost nové sítě	RAM původní	RAM nové
CZ	290 MB	488 MB	26% z 8 GB	34% z 8 GB
RU	299 MB	458 MB	17,6% z 8 GB	24,5% z 8GB

Tabulka 5.4: Vytíženost RAM.

### 5.1.5 Zhodnocení testů

Z testů rozpoznávání je vidět že implementovaný postup generuje rozpoznávací sítě, které jsou větší, pomalejší, a práce s nimi zabírá více operační paměti, ale při jejich použití se dosahuje stejných přesností jako u sítí, které byli vygenerovány skriptem. Důvodem je že implementovaná metoda vynechává poslední krok ze skriptu, kterým je nástroj **SExpand** z projektu STK<sup>3</sup>. Ten implementuje další optimalizační metody.

---

<sup>3</sup>více o STK na webu <http://speech.fit.vutbr.cz/cs/software/hmm-toolkit-stk>

## Kapitola 6

# Rozpoznávač jako modul UniMRCP

K implementaci modulu pro UniMRCP byl použit dekodér pro rozpoznávače pracujícím s aktuálním časem, který pro rozpoznání a přepis textu používá pouze část zvukové stopy od jejího začátku po aktuálně zpracovávaný rámec řeči. Jednoduše si to lze představit jako rozpoznávání řeči pronášené do mikrofону. Tento postup odpovídá zpracováváním krátkých promluv, kde se rozpoznávač musí sám rozhodnout zda rozpoznávání ukončit, či dále pracovat.

### 6.1 Plugin UniMRCP

Plugin je definovaný rozhraním, kterým je struktura jazyka C obsahující ukazatele na funkce<sup>1</sup> a jiné struktury. Rozhraní je definováno pro každý mediální zdroj, tedy rozpoznávání, nahrávání, verifikaci a převod textu na řeč.

UniMRCP obsahuje server a řízení jeho stavu na základě zpráv protokolu MRCP, takže toto nemusí programátor pluginu řešit. Ten se tak může více soustředit na implementaci. Implementovanému se v terminologii UniMRCP říká "engine".

#### 6.1.1 Vytvoření nového pluginu

UniMRCP je postaveno na nástroji `autoconfig`<sup>2</sup>, proto je nutné přidat informace o novém pluginu<sup>3</sup> do konfiguračních souboru v projektu.

#### 6.1.2 Konfigurace pluginu

V souboru `unimrcp/conf/unimrcpserver.xml` je v XML tagu `engine` hlavní konfigurace pluginu. Kde je uloženo jméno enginu, podle něhož se hledá dynamická knihovna s jeho implementací. Tato knihovna se načte a začne pracovat, pokud je parametr `enable` nastaven na `true`. Konfigurační záznam pluginu rozpoznávače pro Phonexii:

<sup>1</sup>Ukazatel na funkci se v terminologii jazyka C také někdy nazývá callback.

<sup>2</sup>Nástroj pro generování make skriptů. Více na webu: <http://www.gnu.org/software/autoconf/>

<sup>3</sup>Podrobnější popis na webu:

<http://www.wikihouse.com/pman10/index.php?Linux%A5%E1%A5%E2%2FUniMRCP%2F%A5%D7%A5%E9%A5%B0%A5%A4%A5%F3%C4%C9%B2%C3%CA%FD%CB%A1>

```

<engine id="Phonexia-ASR" name="phonexiaasr" enable="true">
  <max-channel-count>1</max-channel-count>
  <!-- . is folder where is this xml config file -->
  <param name="bsapi-license-file" value=""/>
  <param name="bsapi-setting-file" value="../../../Phonexia-ASR-data/settings/config_onlinerec_czech"/>
  <param name="bsapi-gc-setting-file" value="../../../Phonexia-ASR-data/settings/config_gc_czech"/>
</engine>

```

V tagu `max-channel-count` je uložen maximální počet VOIP telefonátů, které může engine obsluhovat v jeden okamžik. Další parametry pak definují konfigurační proměnné, které jsou skrze rozhraní pluginu dostupné zevnitř engine.

Přes parametr `bsapi-gc-setting-file` se BSAPI uvnitř pluginu dozví jméno konfiguračního souboru pro proces nastavení pro sestavení rozpoznávací sítě z gramatiky.

### 6.1.3 Rozhraní pluginu rozpoznávače

Obrázek 6.1 zachycuje celé rozhraní, které se dělí na části pro jádro rozpoznávače, jádro modulu a zpracování kanálů.

#### Jádro modulu

Je uloženo ve struktuře `asr_engine_t`. Ukazatele na funkce, které s ním pracují jsou uloženy v `mrcp_engine_method_vtable_t`. Na začátku je zavolána funkce `mrcp_engine_open` ve které se modul inicializuje. Vytvoří tolik rozpoznávacích jader kolik je nastaveno kanálů v konfiguraci.

Funkce `mrcp_engine_channel_create` je volaná, pokaždé když se IVR pokouší o vytvoření nového kanálu. V této funkci je ke kanálu přiřazeno rozpoznávací jádro, které ho obsluhuje.

#### Jádro rozpoznávače

Je uloženo ve struktuře `phonexia_asr_t`. Obsahuje překladač gramatik pro sestavení rozpoznávací sítě a rozpoznávač řeči.

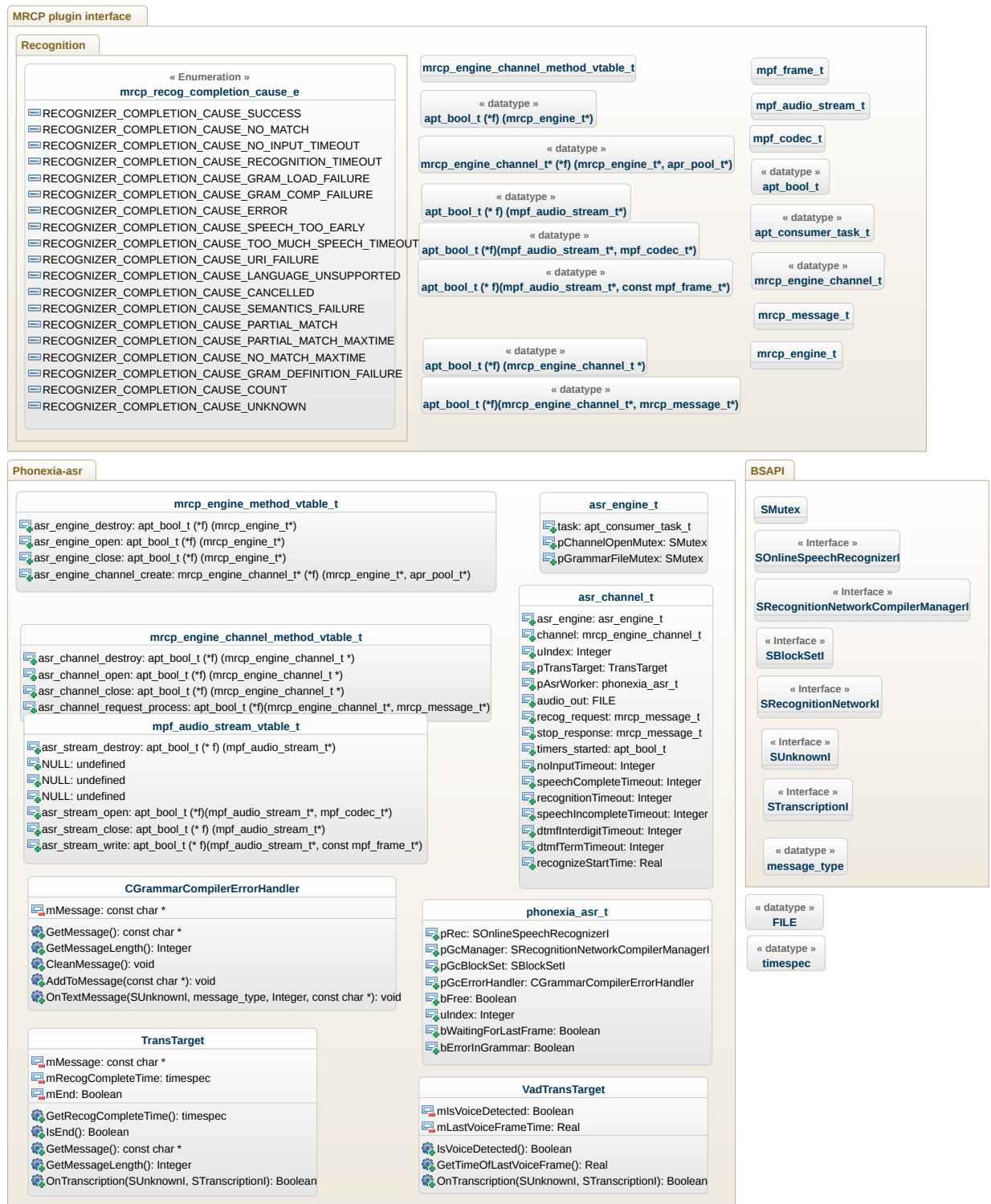
#### Zpracování kanálů

Struktura `asr_channel_t` obsluhuje spojení dialogového manažera s UniMRCP pluginem.

### 6.1.4 Ukončení rozpoznávání

Výčtový typ `mrcp_recog_completion_cause_e` obsahuje definované důvody, proč bylo ukončeno rozpoznávání, nebo komunikace po kanálu.

Plugin díky kombinací gramatiky a nastavených časových limitů, sám určí kdy je rozpoznávání ukončeno a pošle výsledek dialogovému manažerovi.



Obrázek 6.1: Diagram UniMRCP modulu pro ASR.

## Kapitola 7

# Použití UniMRCP ve FreeSwitch

O komunikaci UniMRCP a FreeSwitch se stará modul `mod.unimrcp`<sup>1</sup>, který je součástí projektu FreeSwitch.

Pro použití rozpoznávače řeči je zapotřebí nakonfigurování spojení s MRCP serverem, vytvoření dialogových scénářů a napojení tohoto scénáře na SIP telefonní číslo. Dialogový scénář je aplikací, kde vzniká největší přidaná hodnota pro uživatele, protože to on definuje jak systém zareaguje na výsledek rozpoznávání řeči.

### 7.1 Konfigurace spojení s UniMRCP

Se musí definovat v XML souboru ve složce `conf/mrcp_profiles`. Obsah souboru `phonexiamrcp.xml` pro spojení s UniMRCP:

```
<include>
  <profile name="phonexiamrcp" version="2">
    <param name="client-ip" value="127.0.0.1"/>
    <param name="client-port" value="5090"/>
    <param name="server-ip" value="127.0.0.1"/>
    <param name="server-port" value="8060"/>
    <param name="sip-transport" value="udp"/>
    <param name="ua-name" value="FreeSWITCH"/>
    <param name="rtp-ip" value="127.0.0.1"/>
    <param name="rtp-port-min" value="4000"/>
    <param name="rtp-port-max" value="5000"/>
    <param name="codecs" value="PCMU_LPCMA_L16/96/8000"/>
    <recogparams>
    </recogparams>
  </profile>
</include>
```

Klient v konfiguraci znamená MRCP klient, tedy ten kdo využívá služeb MRCP serveru. Tím je FreeSwitch na kterém je tato konfigurace. Jednotlivé parametry popisuje tabulka 7.1.

<sup>1</sup>Více na webu: [https://wiki.freewitch.org/wiki/Mod\\_unimrcp](https://wiki.freewitch.org/wiki/Mod_unimrcp)

<sup>2</sup>IP (Internet Protocol) adresa slouží pro adresování počítače v internetové síti. Více v RFC791 na webu: <http://tools.ietf.org/html/rfc791>



Parametr	Význam
client-ip	Ip <sup>2</sup> adresa FreeSwitch serveru.
client-port	Port FreeSwitch serveru pro protokol MRCP.
server-ip	Ip adresa UniMRCP serveru.
server-port	Port UniMRCP serveru pro protokol MRCP.
sip-transport	protokol pro přenos SIPu. UDP <sup>3</sup> nebo TCP <sup>4</sup> .
codecs	Seznam formátů pro přenos audia.
rtp-port-min rtp-port-max	Rozsah portů na serveru UniMRCP pro přenos audia.

Tabulka 7.1: Tabulka parametrů pro konfiguraci MRCP spojení.

## 7.2 Dialogový scénář

Lze v projektu FreeSwitch sestavit v Perlu, jazyku C, Lua skriptu a Javascriptu. Scénáře jsou uloženy ve složce freeswitch/scripts. Scénář pro volání ASR definovaný v Lua skriptu vypadá následovně:

```
session:answer()
session:execute("play_and_detect_speech","say:_detect:unimrcp_
    {no-input-timeout=5000_,recognition-timeout=5000_,start-
    input-timers=true,define-grammar=true}_test" )

local speech_result = session:getVariable('
    detect_speech_result')
if (result ~= nil) then
    freeswitch.consoleLog("INFO", speech_result .. "\n")
    message:chat_execute("reply", speech_result);
else
    freeswitch.consoleLog("INFO","No_result!\n")
end
```

Na začátku je příkazem `session:answer()` zvednuto sluchátko v ústředně. Skript spustí pomocí funkce `execute` spustí metodu `play_and_detect_speech`, který pomocí TTS přehraje argument `say`. V tomto volání není vyplněn. A po přehrání požádá server MRCP o rozpoznání řeči podle gramatiky `test`. Pokud došlo k úspěšnému rozpoznání je v proměnné `detect_speech_result` uložen výsledek vrácený z MRCP, který se načte do proměnné `speech_result`. V opačném případě se nic nenačte a v proměnné je `nil`.

## 7.3 Přiřazení dialogového scénáře k telefonnímu číslu

Je definováno v XML souboru `freeswitch/conf/dialplan/public.xml` následovně:

```
<extension name="Phonexia_ASR">
  <condition field="destination_number" expression="^4001$" >
```

<sup>3</sup>Protokol pro přenos dat UDP (User Datagram protokol). Více v RFC768 na webu: <http://tools.ietf.org/html/rfc768>.

<sup>4</sup>Protokol pro přenos dat TCP (Transmission Control Protocol). Více v RFC793 na webu: <http://tools.ietf.org/html/rfc793>.

```
<action application="lua" data="asrdevel.lua" />
</condition>
</extension>
```

Záznam říká, že při volání čísla 4001 se má spustit dialogový scénář v souboru `asrdevel.lua`, který je napsaný v Lua skriptu. FreeSwitch vyhledává scénáře v Lua skriptu ve složce `freeswitch/scripts`.

## Kapitola 8

# Závěr

Výsledkem této práce je základ IVR systému pro rozpoznávač řeči řízený SRGS/ABNF gramatikou. Typickým použitím tohoto systému je nahrazení odpovědi pomocí DTMF volby použitím lidské řeči. Dalším použitím je realizace hlasového vstupu pro zařízení ovládána hlasem.

Pro realizování tohoto systému jsem vytvořil překladač podmnožiny SRGS/ABNF gramatik na váhovaný konečný automat. Projekt BSAPI jsem doplnil, tak aby podporoval překladač gramatik. Konečný automat je v něm podporován knihovnou OpenFst. Podle existujících skriptů a článků [7] jsem implementoval sestavování rozpoznávacích sítí v C++ v BSAPI. S použitím BSAPI dekodéru jsem vytvořil plugin pro rozpoznávání řeči pro UniMRCP projekt. Plugin komunikuje se softwarovou telefonní ústřednou FreeSwitch pomocí protokolů MRCPv2 a RTP. Ve FreeSwitchi jsem nakonfiguroval spojení s pluginem pro rozpoznávání řeči a definoval jsem dialogový scénář pro prezentaci systému.

V překladači jsem neimplementoval podporu některých elementů jazyka třeba podpora vícejazykových gramatik, ale především elementu "tag", který slouží pro definování sémantické interpretace řeči. Překladač ani neakceptuje vstup v XML formátu. O tuto funkcionalitu by měl být v překladači v budoucnu rozšířen. Dále doplnit implementaci pluginu pro rozpoznávání řeči, jelikož plně nepodporuje zcela všechny možnosti protokolu MRCPv2, a to především důvěryhodnost s jakou rozpoznávač text určil. Proto ale bude potřeba o tuto funkci rozšířit dekodér pro online rozpoznávání řeči v BSAPI. A v neposlední řadě zahrnout optimalizaci z nástroje **SExpand** z projektu STK do procesu sestavování sítí.

Díky tomuto projektu jsem získal nové zkušenosti z oboru zpracování přirozeného jazyka, konfigurace IP telefonie. Prohloubil své znalosti z oblastí teorie konečných automatů a síťových protokolů. Ale především jsem se seznámil s prací ve firemním prostředí, kde jsem mohl sledovat zpětnou vazbu manažerů a zákazníků na mou práci.

# Literatura

- [1] A. DEROUAULT, B. M.: Natural Language Modeling for Phoneme-to-Text Transcription.  
URL <http://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=04767855>
- [2] A. Hunt, S. McGlashan: Speech Recognition Grammar Specification Version 1.0.  
URL <http://www.w3.org/TR/speech-grammar/>
- [3] D. Burnett, S. S.: Media Resource Control Protocol Version 2 (MRCPv2), 2012.  
URL <http://tools.ietf.org/html/rfc6787>
- [4] D. Povey: The Kaldi Speech Recognition Toolkit.  
URL <http://homepages.inf.ed.ac.uk/aghoshal/pubs/asru11-kaldi.pdf>
- [5] Freeform Dynamics: 2005 Speech Recognition Customer Satisfaction Survey. Technická zpráva, 2005.  
URL [http://www.freeformdynamics.com/misc/docs/Nuance\\_Customer\\_Sat\\_ExecSumm.pdf](http://www.freeformdynamics.com/misc/docs/Nuance_Customer_Sat_ExecSumm.pdf)
- [6] M. Mohri: Weighted Finite-State Transducer Algorithms.  
URL <http://www.cs.nyu.edu/~mohri/pub/fla.pdf>
- [7] M. Mohri, F. Pereira, M. R.: SPEECH RECOGNITION WITH WEIGHTED FINITE-STATE TRANSDUCERS.  
URL <http://www.cs.nyu.edu/~mohri/pub/hbka.pdf>
- [8] R. KUHN, R. M.: A Cache-Based Natural Language Model for Speech Recognition.  
URL <http://ieeexplore.ieee.org/stamp/stamp.jsp?tp=&arnumber=56193>
- [9] S. J. Young, L. L. C.: Speech recognition evaluation: a review of the U.S. CSR and LVCSR programmes.  
URL [http://ac.elsa-cdn.com/S0885230898901012/1-s2.0-S0885230898901012-main.pdf?\\_tid=0388f02c-e4ec-11e3-91cc-00000aab0f27&acdnat=1401118987\\_0c0bb8b00af46f7872d95d660e56f2ac](http://ac.elsa-cdn.com/S0885230898901012/1-s2.0-S0885230898901012-main.pdf?_tid=0388f02c-e4ec-11e3-91cc-00000aab0f27&acdnat=1401118987_0c0bb8b00af46f7872d95d660e56f2ac)
- [10] T. Parr: *The Definitive ANTLR Reference: Building Domain-Specific Languages*. Pragmatic Programmers, Pragmatic Bookshelf, první vydání, Květen 2007, ISBN 0978739256.  
URL <http://www.amazon.com/Definitive-ANTLR-Reference-Domain-Specific-Programmers/dp/0978739256%3FSubscriptionId%3D13CT5CVB80YFWJEPWS02%26tag%3Dws%26linkCode%3Dxm2%26camp%3D2025%26creative%3D165953%26creativeASIN%3D0978739256>

## Příloha A

### Obsah CD

- Soubor "CTIME.txt".
- Zdrojové kódy pro převod SRGS/ABNF na konečný automat.
- Zdrojové kódy části BSAPI, které byly upraveny v rámci projektu.  
BSAPI nelze spustit bez licence, pro bližší informace kontaktujte [support@phonexia.com](mailto:support@phonexia.com)
- Zdrojové kódy modulu rozpoznávače řeč pro UniMRCP.